

Regression

Classification

- Input: X
 - Real valued, vectors over real.
 - Discrete values (0,1,2,...)
 - Other structures (e.g., strings, graphs, etc.)
- Output: Y
 - Discrete (0,1,2,...)

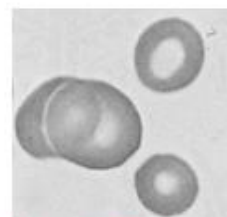
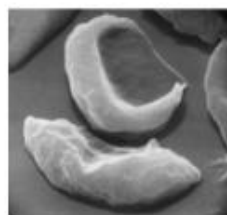


$X = \text{Document}$



Sports
Science
News

$Y = \text{Topic}$



$X = \text{Cell Image}$

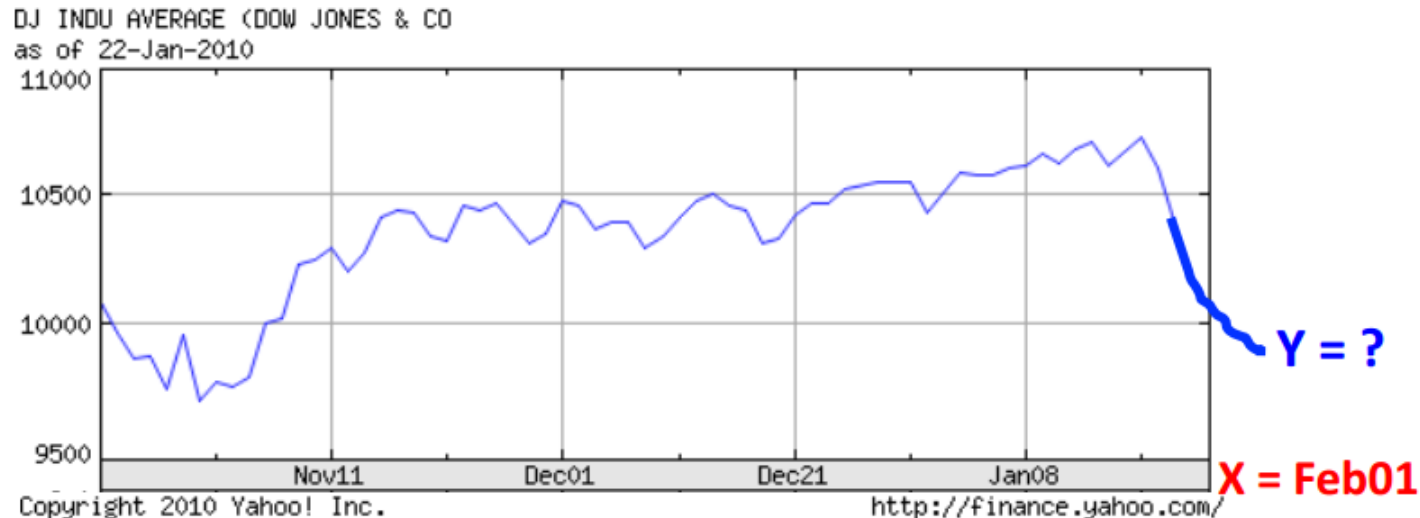


Anemic cell
Healthy cell

$Y = \text{Diagnosis}$

- Input: X
 - Real valued, vectors over real.
 - Discrete values (0,1,2,...)
 - Other structures (e.g., strings, graphs, etc.)
- Output: Y
 - Real valued, vectors over real.

Stock Market
Prediction



What should I watch tonight?

All

[Movies, TV & Showtimes](#) [Celebs, Events & Photos](#) [News & Community](#) [Watchlist](#)



Point Break (2015)

PG-13 | 114 min | 25 December 2015

5.4 **Our rating:** ★★★★★★ ★★ -/10

Ratings: **5.4/10** from 7,322 users Metascore: 34/100

Reviews: 60 user | 84 critic | 19 from Metacritic.com

A young FBI agent infiltrates an extraordinary team of extreme sports athletes he suspects of masterminding a string of unprecedented, sophisticated corporate heists. "Point Break" is inspired by the classic 1991 hit.

Director: [Ericson Core](#)

Writers: [Kurt Wimmer](#) (screenplay), [Rick King](#) (story), [5 more credits](#) »

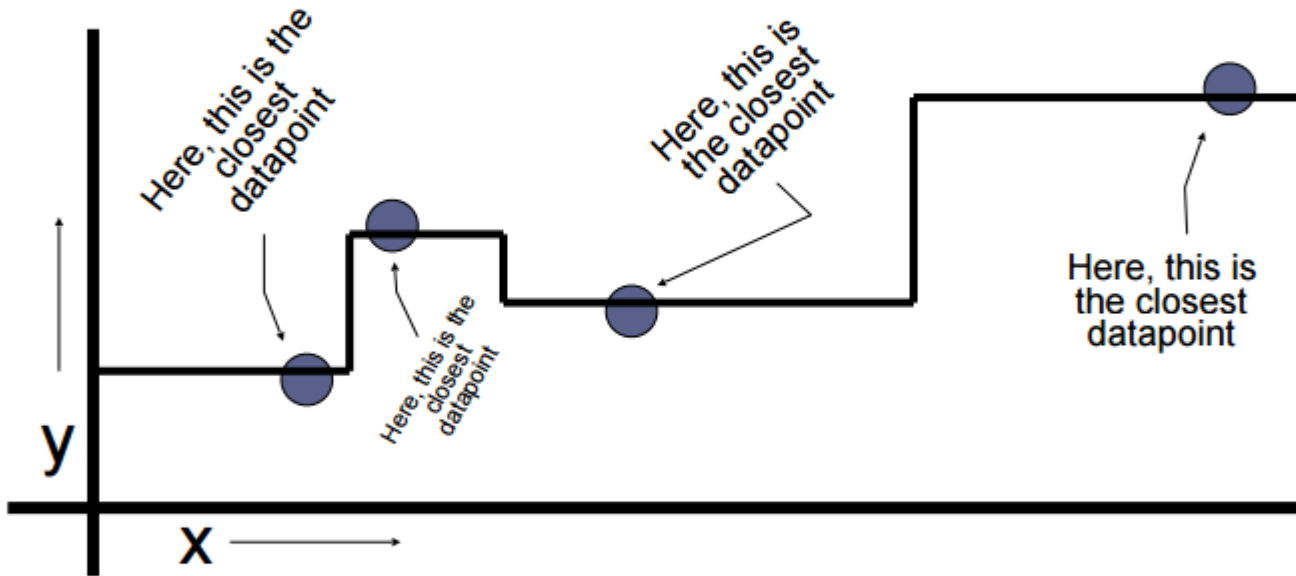
Stars: [Édgar Ramírez](#), [Luke Bracey](#), [Ray Winstone](#) | [See full cast and crew](#) »

[+ Watchlist](#) [Watch Trailer](#) [Share...](#)

[See More on IMDb Pro](#) »

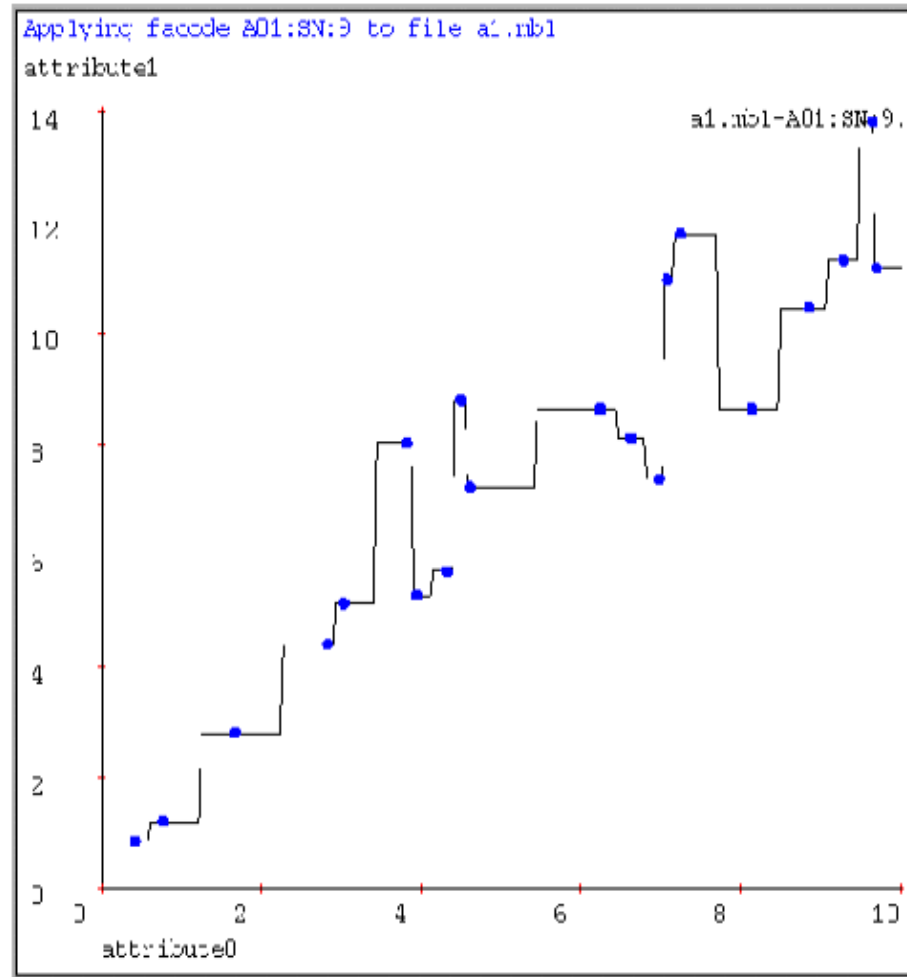
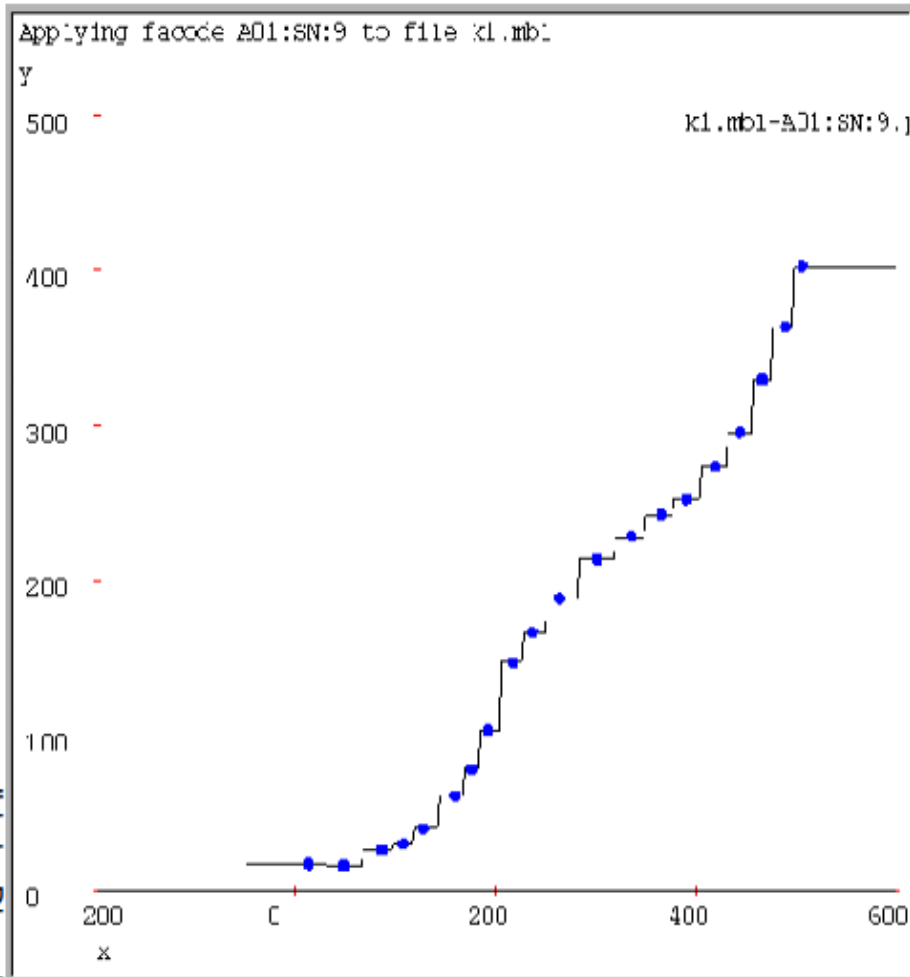
Predict this automatically!

1-NN for Regression



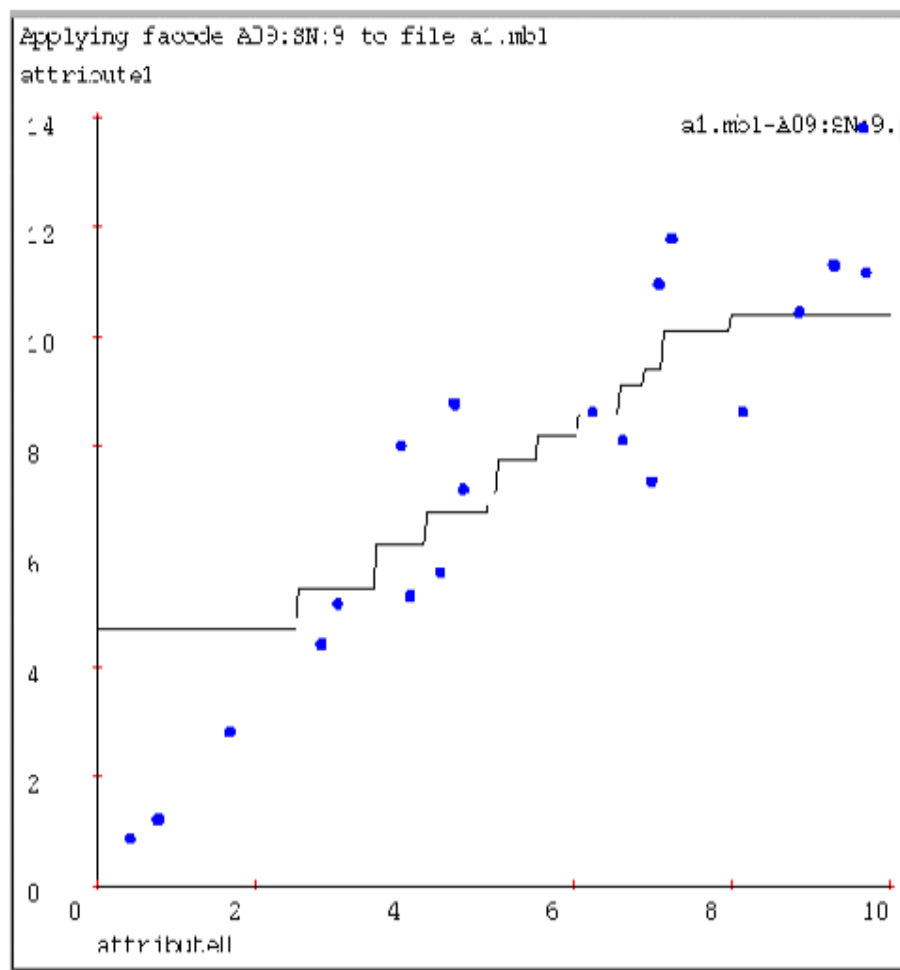
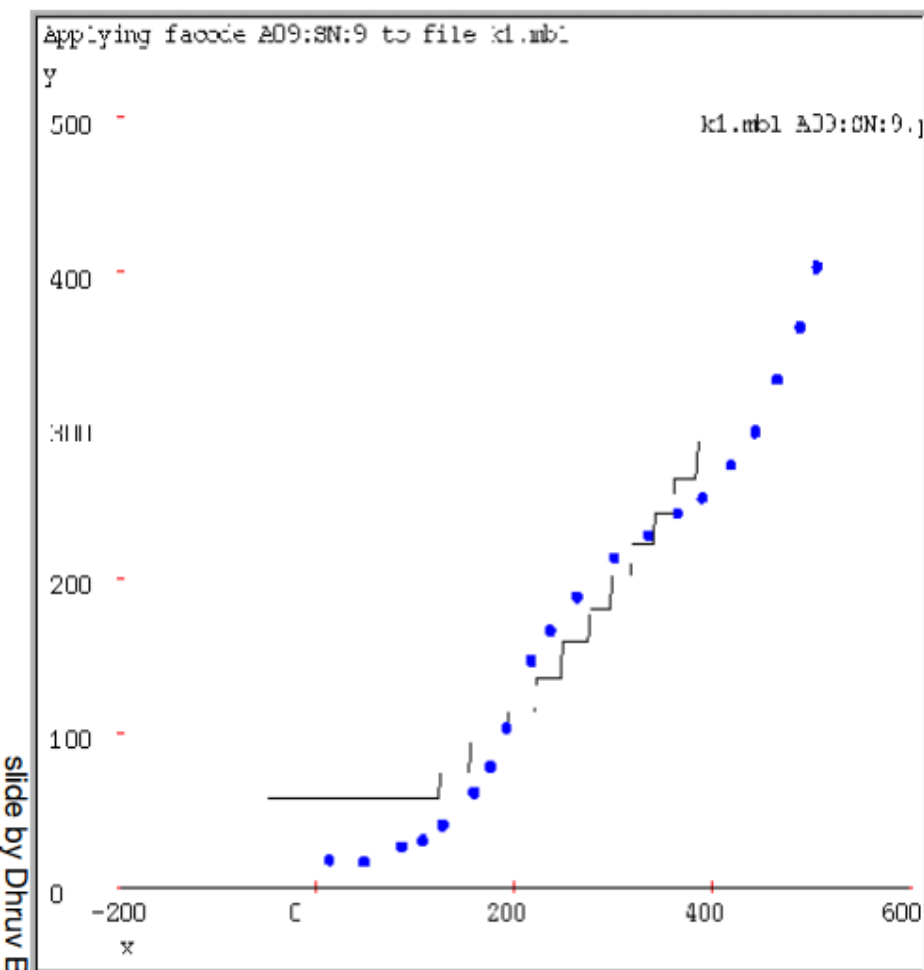
1-NN for Regression

- Often bumpy (overfits)

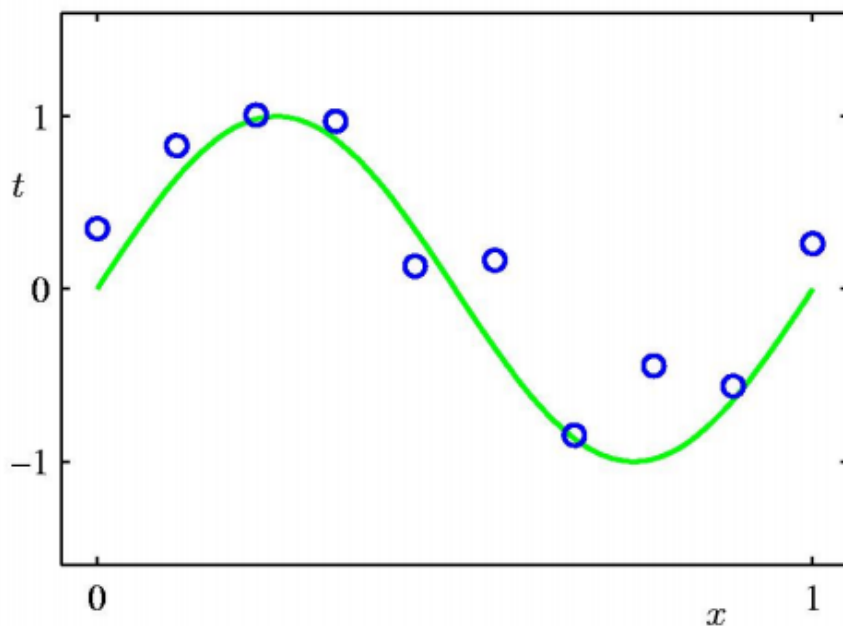


9-NN for Regression

- Often bumpy (overfits)



Simple 1-D Regression



- Circles are data points (i.e., training examples) that are given to us
- The data points are uniform in x , but may be displaced in y

$$t(x) = f(x) + \varepsilon$$

with ε some noise

- In **green** is the “true” curve that we don’t know

What is a Model?

1. Often Describe Relationship between Variables

2. Types

- Deterministic Models (no randomness)
- Probabilistic Models (with randomness)

Deterministic Models

1. Hypothesize Exact Relationships
2. Suitable When Prediction Error is Negligible
3. Example: Body mass index (BMI) is measure of body fat based

- $BMI = \frac{\text{Weight in Kilograms}}{(\text{Height in Meters})^2}$

Probabilistic Models

1. Hypothesize 2 Components
 - Deterministic
 - Random Error
2. Example: Systolic blood pressure of newborns Is 6 Times the Age in days + Random Error
 - $SBP = 6 \times \text{age}(d) + \varepsilon$
 - Random Error May Be Due to Factors Other Than age in days (e.g. Birthweight)

Simple Regression

- Simple regression analysis is a statistical tool that gives us the ability to estimate the mathematical relationship between a dependent variable (usually called y) and an independent variable (usually called x).
- The dependent variable is the variable for which we want to make a prediction.
- While various non-linear forms may be used, simple linear regression models are the most common.

Introduction

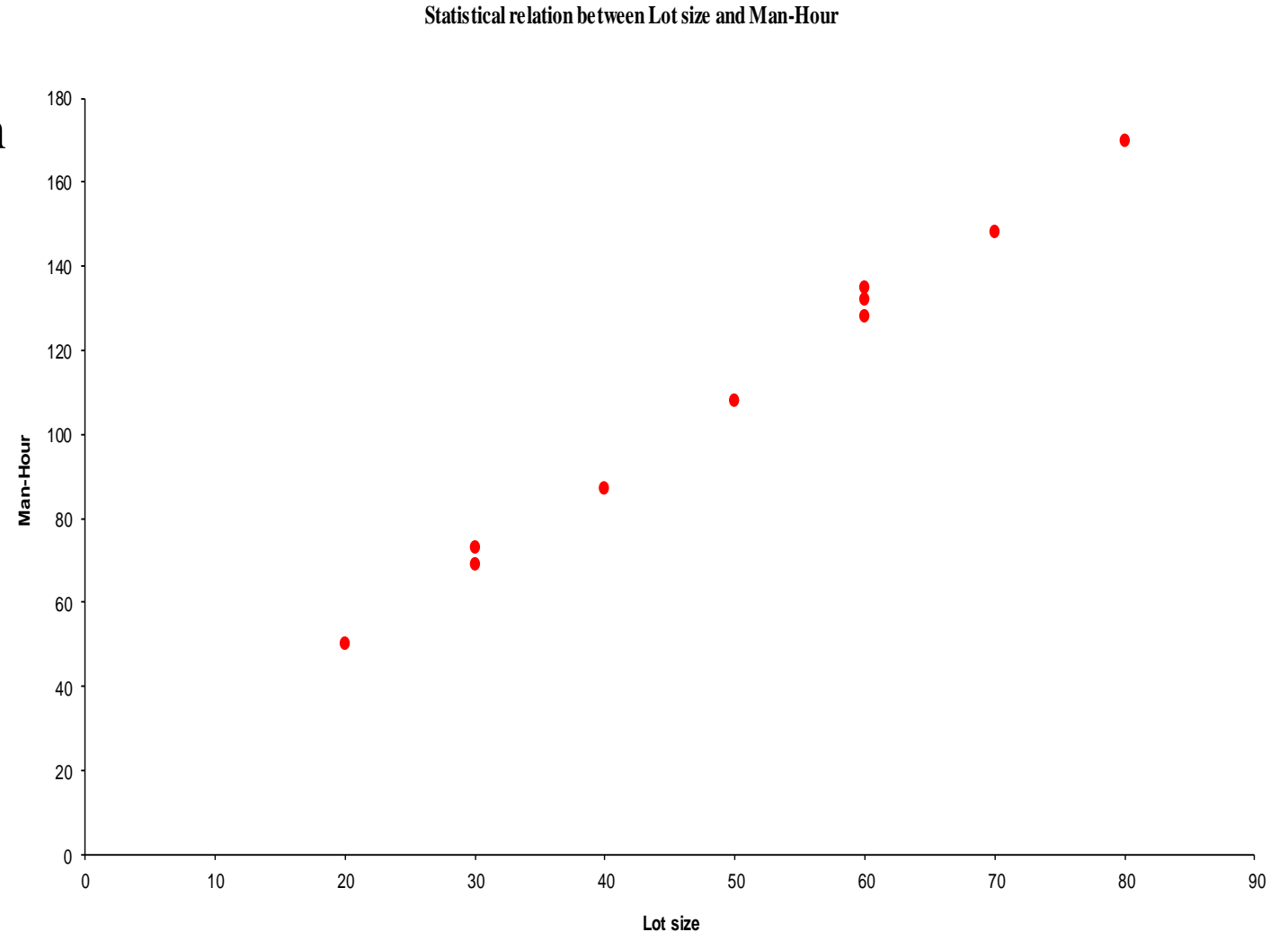
- The primary goal of quantitative analysis is to use current information about a phenomenon to predict its future behavior.
- Current information is usually in the form of a set of data.
- In a simple case, when the data form a set of pairs of numbers, we may interpret them as representing the observed values of an independent (or predictor or explanatory) variable X and a dependent (or response or outcome) variable Y .

lot size	Man-hours
30	73
20	50
60	128
80	170
40	87
50	108
60	135
30	69
70	148
60	132

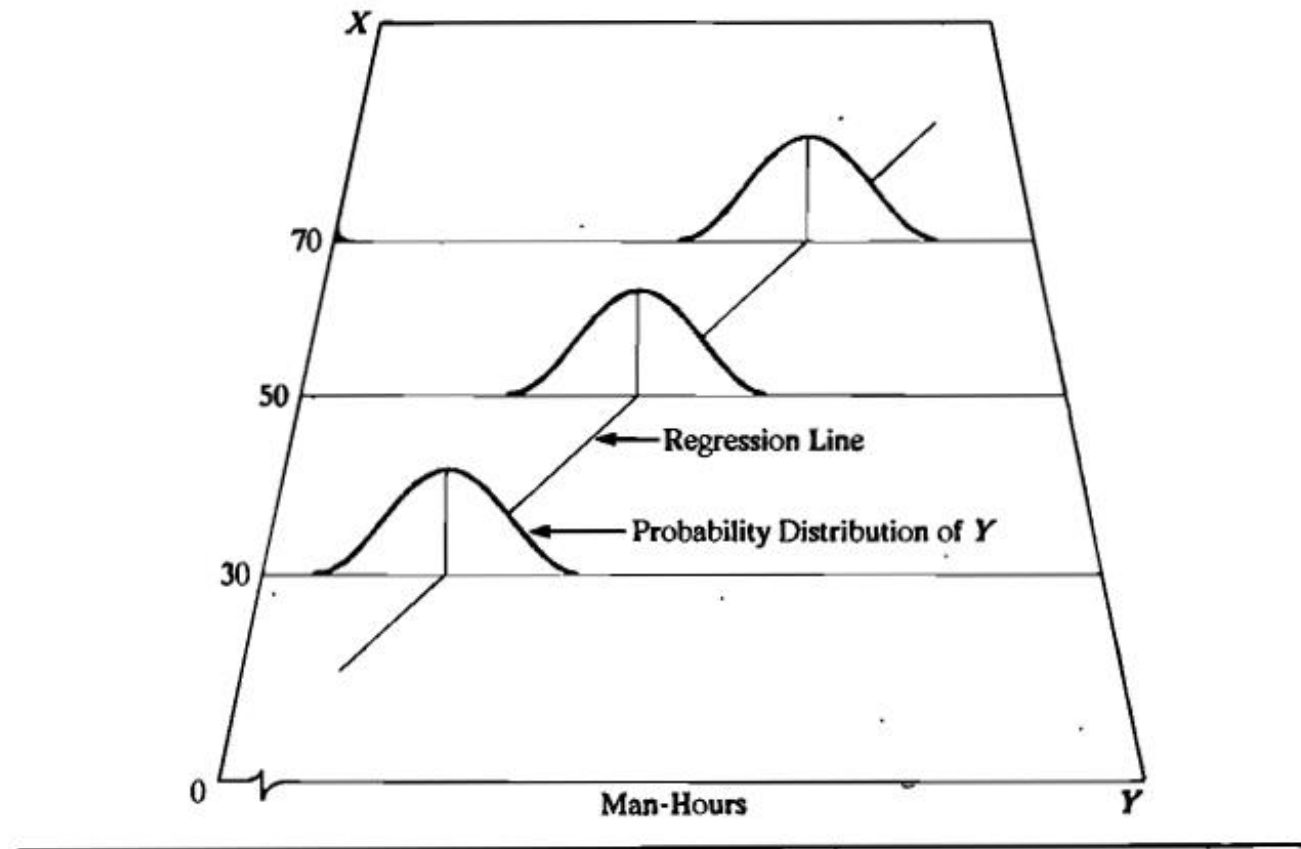
Introduction

- The goal of the analyst who studies the data is to find a functional relation between the response variable y and the predictor variable x .

$$y = f(x)$$

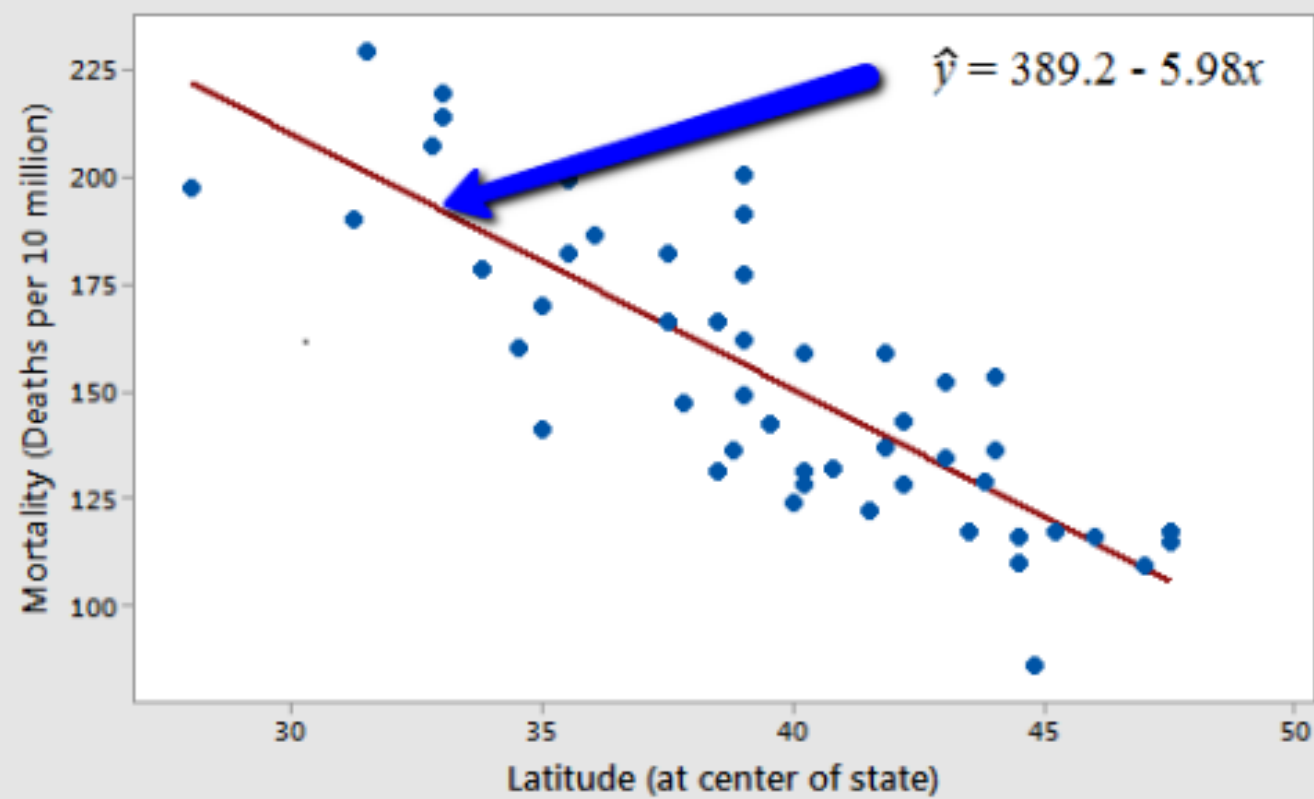


Pictorial Presentation of Linear Regression Model



Linear Regression Model

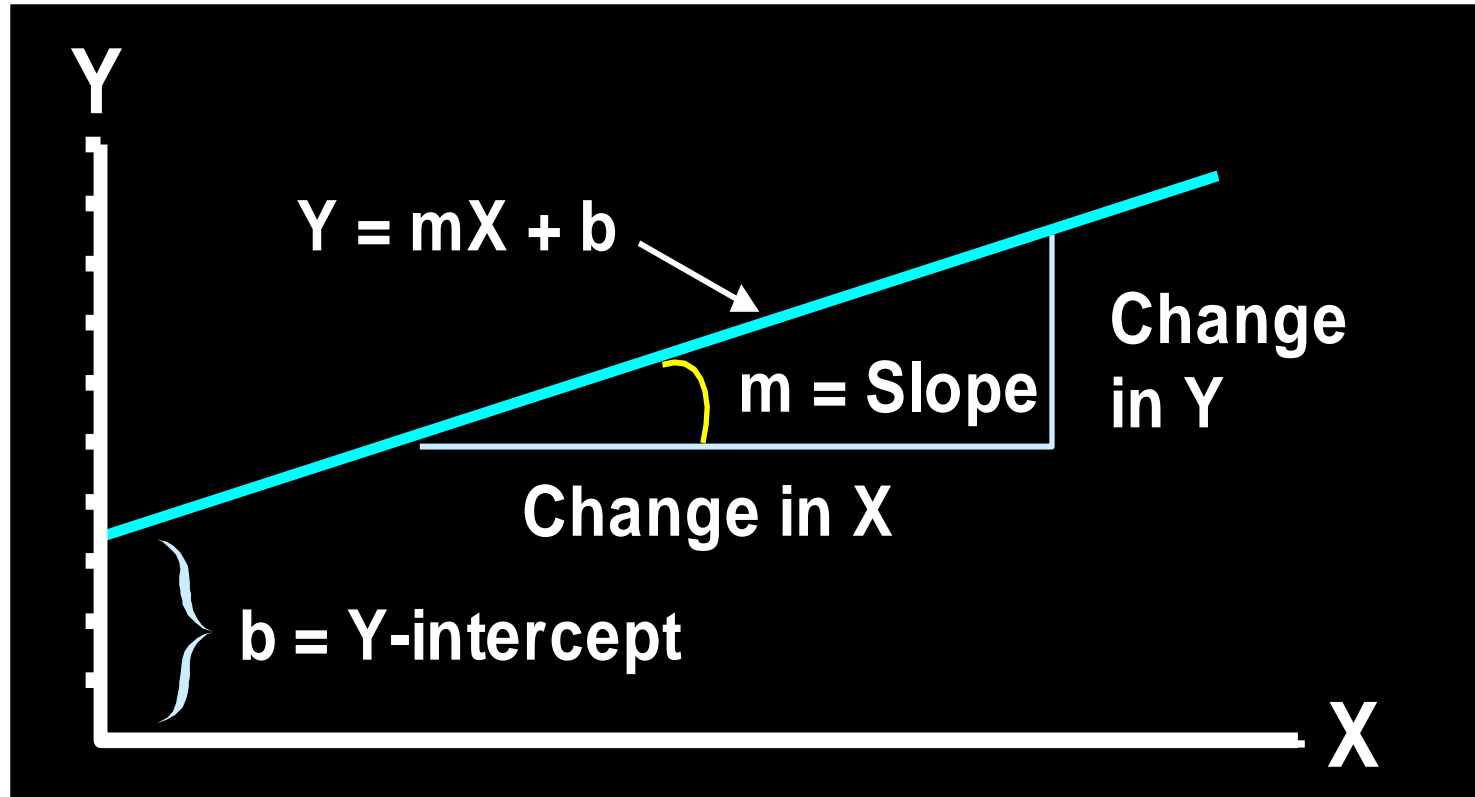
Skin Cancer Mortality versus State Latitude



Assumptions

- Linear regression assumes that...
 - 1. The relationship between X and Y is linear
 - 2. Y is distributed normally at each value of X
 - 3. The variance of Y at every value of X is the same (homogeneity of variances)
 - 4. The observations are independent

Linear Equations



Linear Regression Model

- 1. Relationship Between Variables Is a Linear Function

The diagram illustrates the Linear Regression Model equation $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$. It features five labels with arrows pointing to the corresponding parts of the equation:

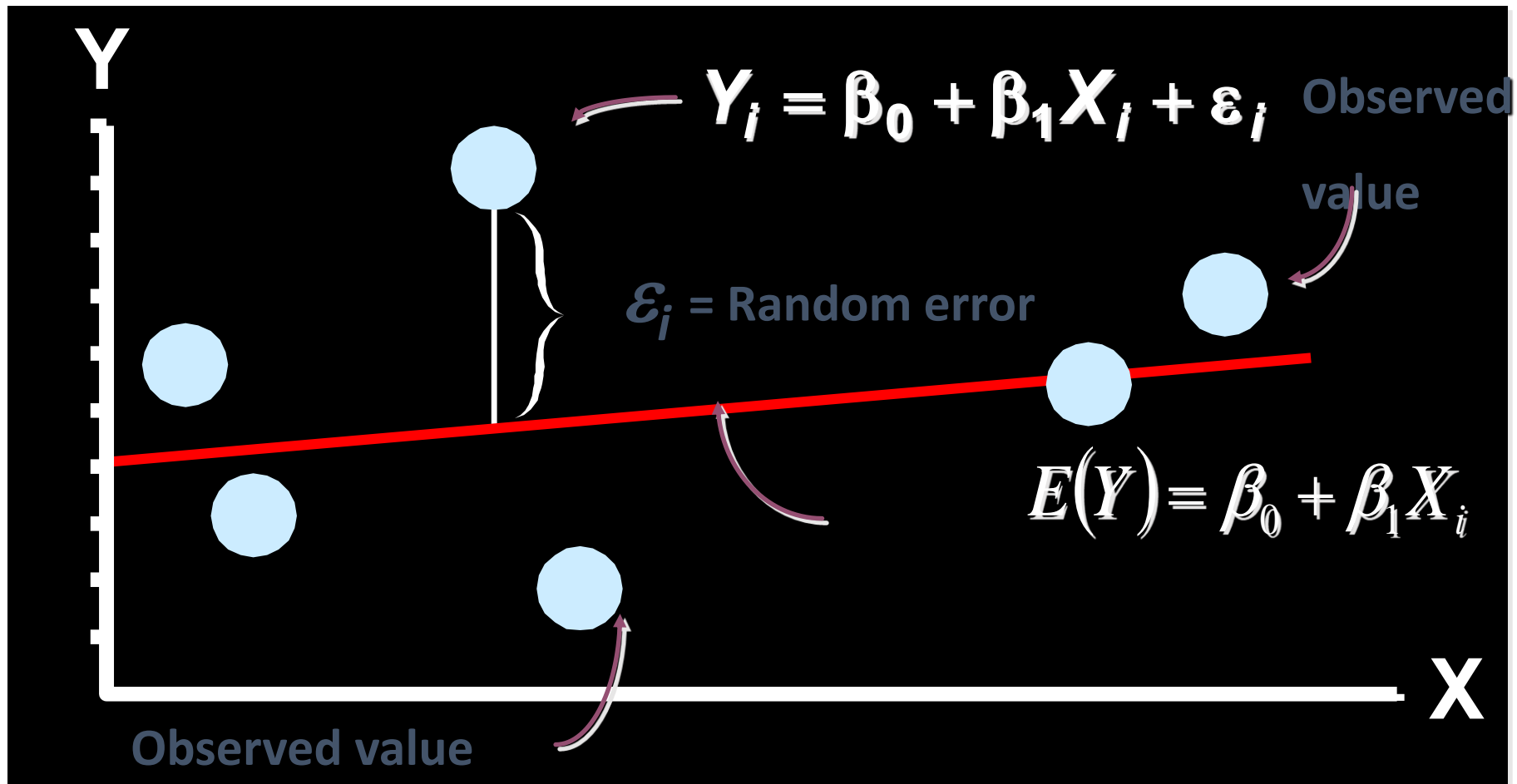
- Population Y-Intercept** points to β_0 .
- Population Slope** points to β_1 .
- Random Error** points to ε_i .
- Dependent (Response) Variable (e.g., CD+ c.)** points to Y_i .
- Independent (Explanatory) Variable (e.g., Years s. serocon.)** points to X_i .

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

Meaning of Regression Coefficients

- General regression model
 1. β_0 , and β_1 are parameters
 2. X is a known constant
 3. Deviations ε are independent $N(0, \sigma^2)$
- The values of the regression parameters β_0 , and β_1 are not known. We estimate them from data.
- β_1 indicates the change in the mean response per unit increase in X .

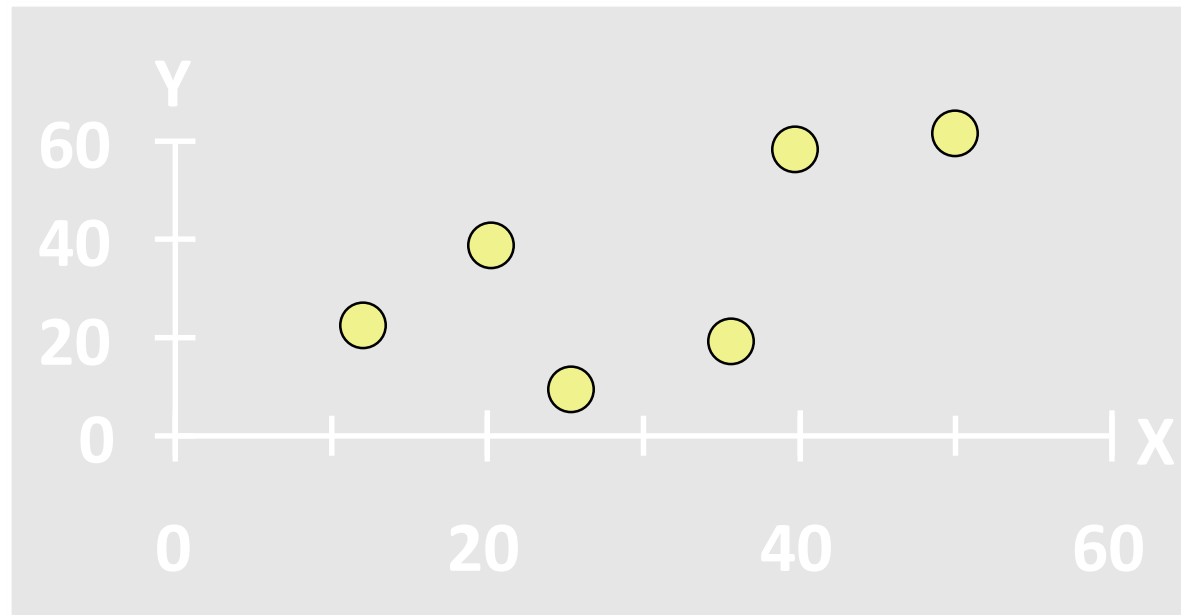
Population Linear Regression Model



Estimating Parameters: Least Squares Method

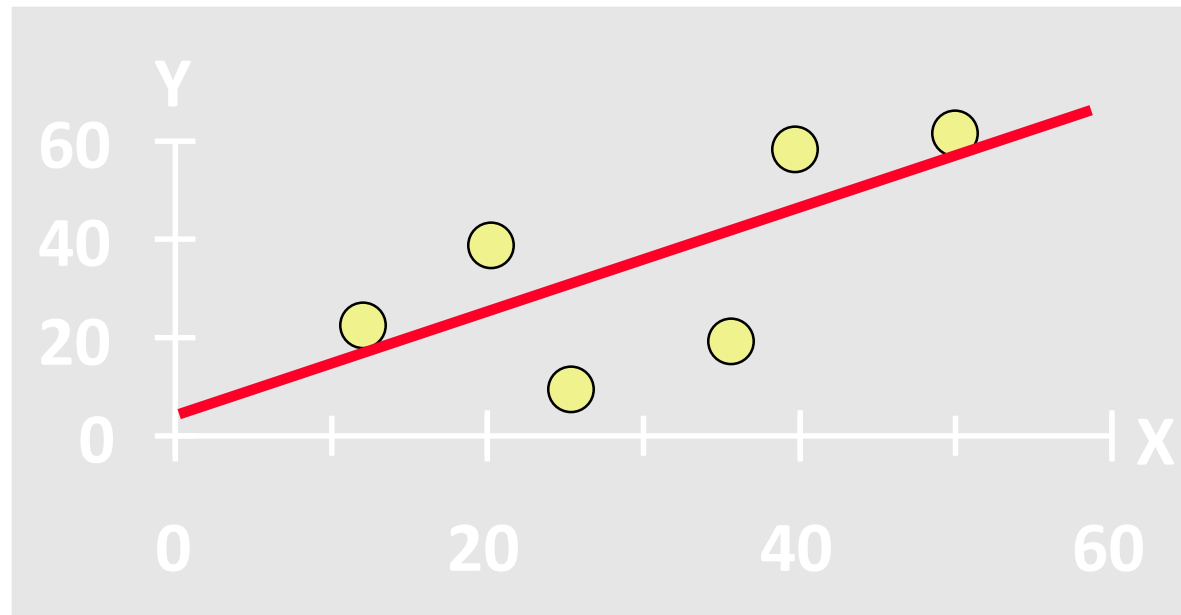
Scatter plot

- 1. Plot of All (X_i, Y_i) Pairs
- 2. Suggests How Well Model Will Fit



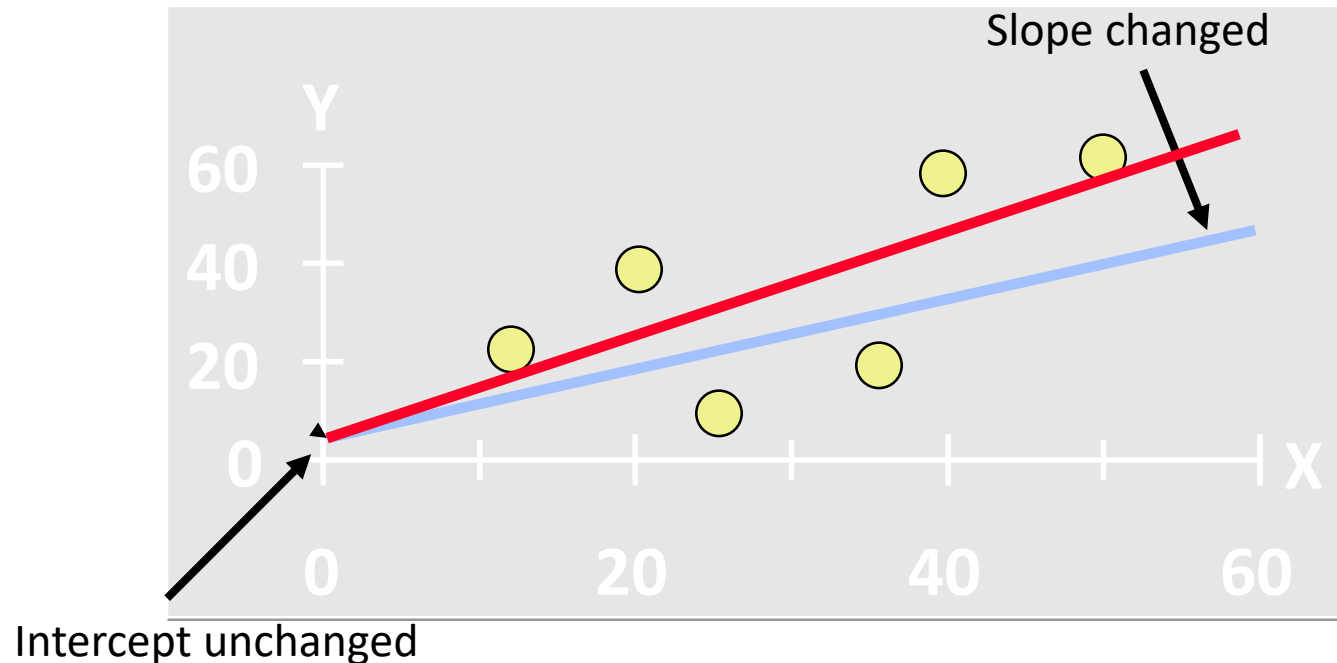
Thinking Challenge

**How would you draw a line through the points?
How do you determine which line 'fits best'?**



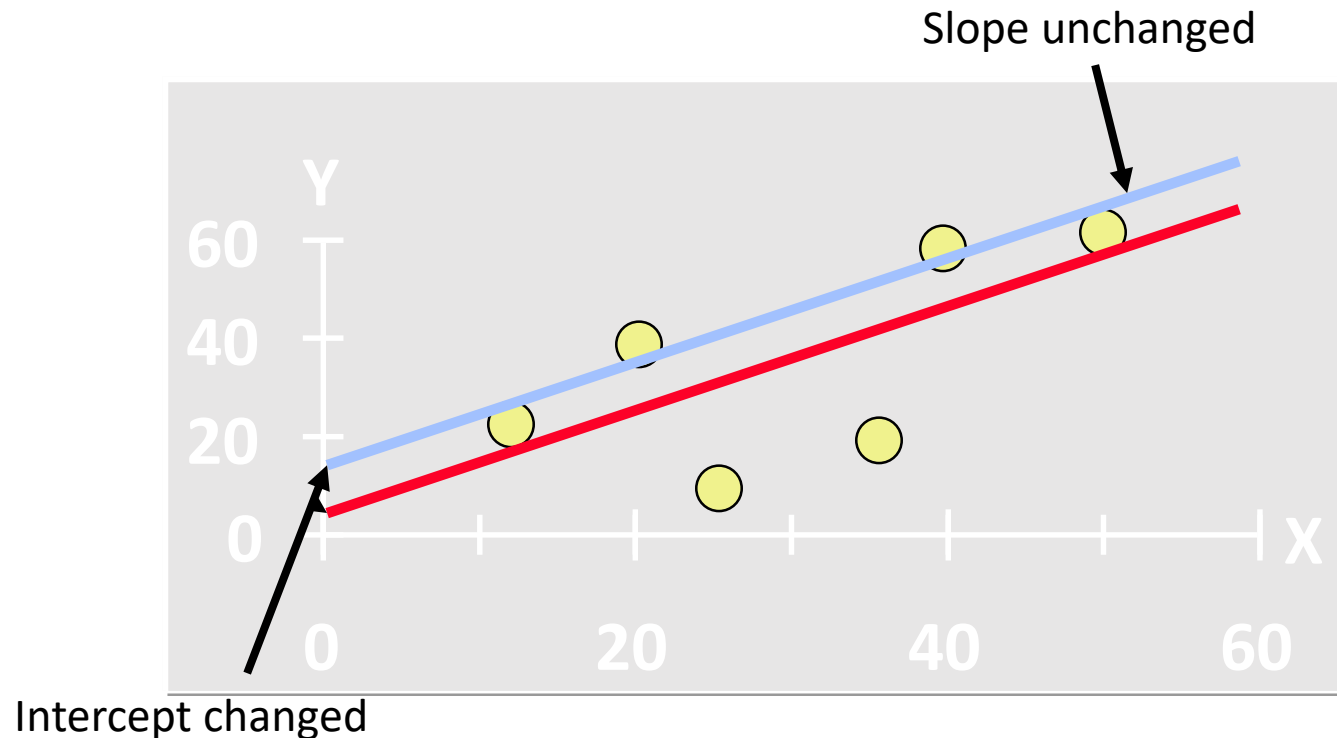
Thinking Challenge

**How would you draw a line through the points?
How do you determine which line 'fits best'?**



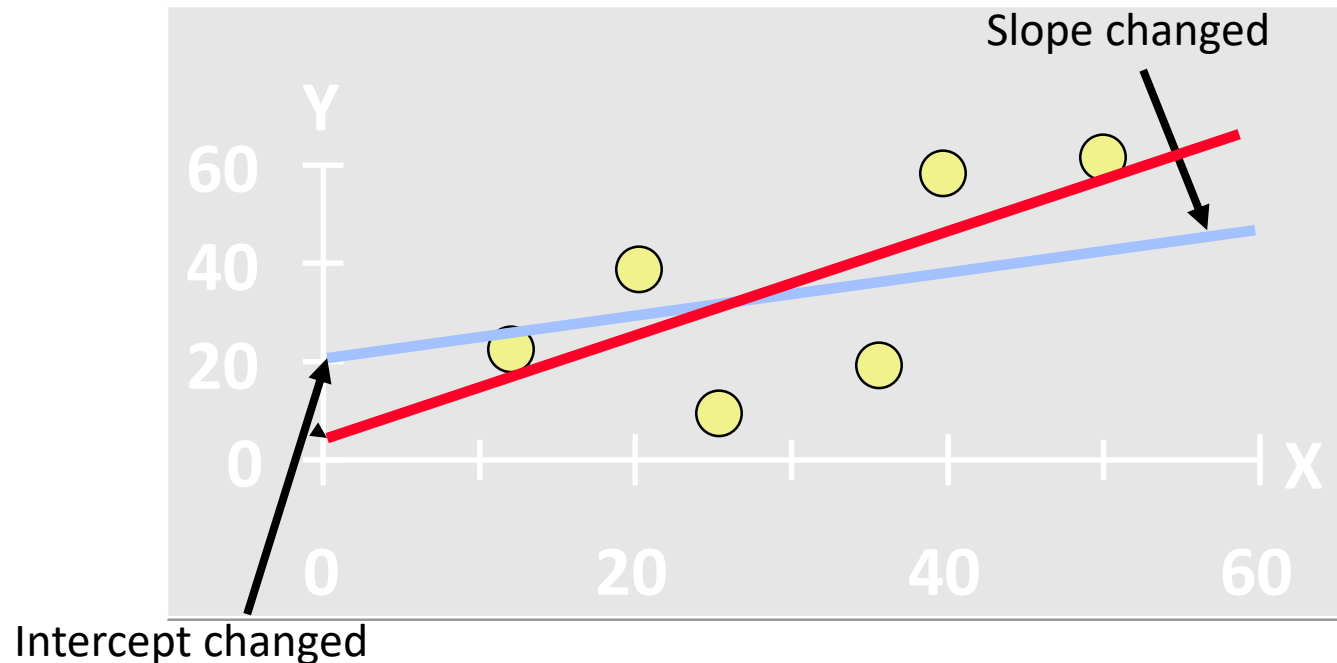
Thinking Challenge

**How would you draw a line through the points?
How do you determine which line 'fits best'?**



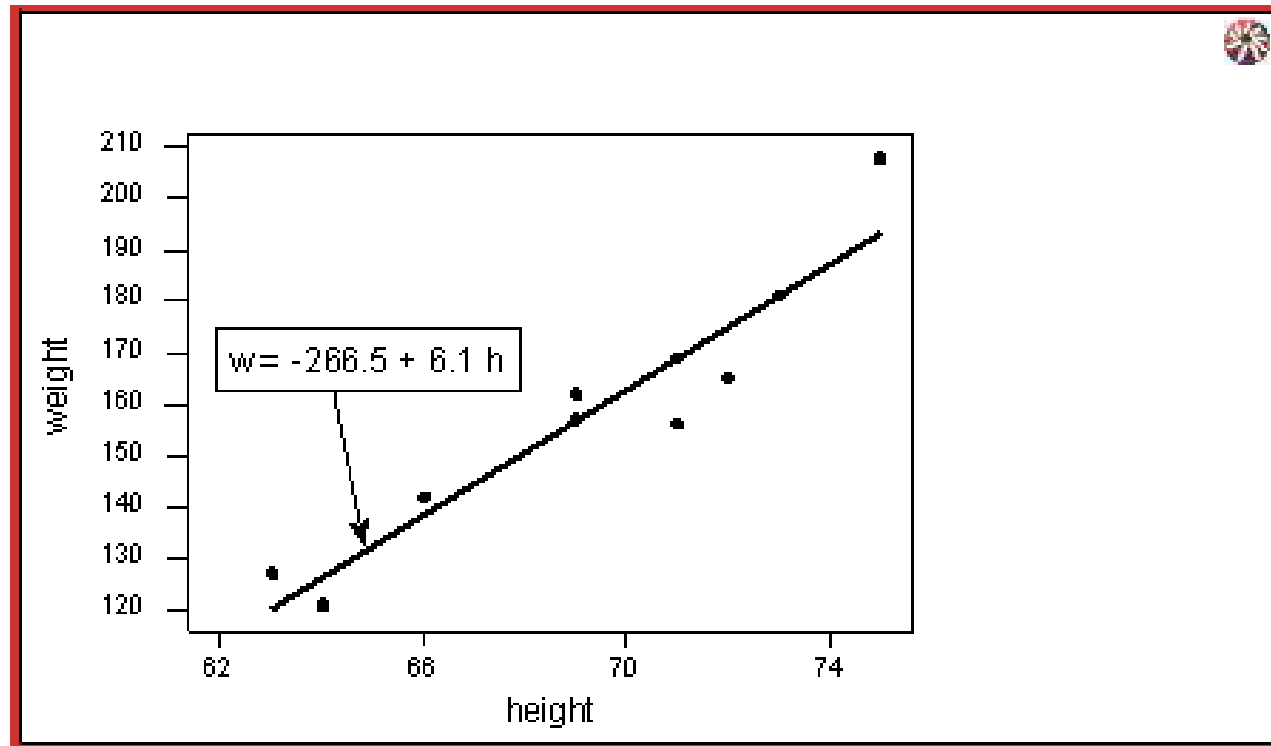
Thinking Challenge

How would you draw a line through the points?
How do you determine which line 'fits best'?



What is the best fitting line

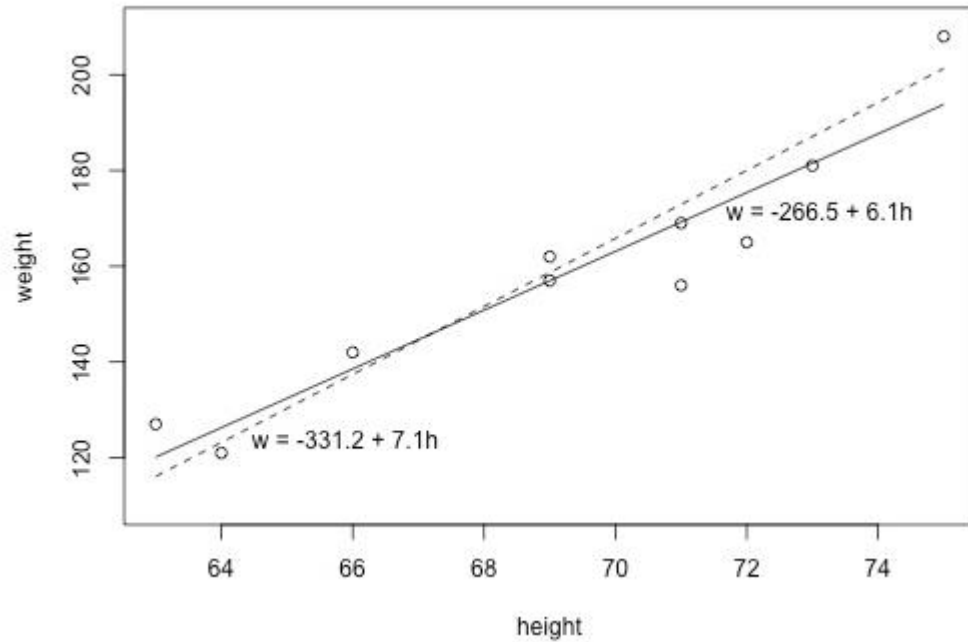
i	x_i	y_i	\hat{y}_i
1	63	127	120.1
2	64	121	126.3
3	66	142	138.5
4	69	157	157.0
5	69	162	157.0
6	71	156	169.2
7	71	169	169.2
8	72	165	175.4
9	73	181	181.5
10	75	208	193.8



$$\hat{y}_i = b_0 + b_1 x_i$$

- y_i denotes the observed response for experimental unit i
- x_i denotes the predictor value for experimental unit i
- \hat{y}_i is the predicted response (or fitted value) for experimental unit i

Prediction Error



$w = -331.2 + 7.1 h$ (the dashed line)					
i	x_i	y_i	\hat{y}_i	$(y_i - \hat{y}_i)$	$(y_i - \hat{y}_i)^2$
1	63	127	116.1	10.9	118.81
2	64	121	123.2	-2.2	4.84
3	66	142	137.4	4.6	21.16
4	69	157	158.7	-1.7	2.89
5	69	162	158.7	3.3	10.89
6	71	156	172.9	-16.9	285.61
7	71	169	172.9	-3.9	15.21
8	72	165	180.0	-15.0	225.00
9	73	181	187.1	-6.1	37.21
10	75	208	201.3	6.7	44.89
					766.5

$w = -266.53 + 6.1376 h$ (the solid line)					
i	x_i	y_i	\hat{y}_i	$(y_i - \hat{y}_i)$	$(y_i - \hat{y}_i)^2$
1	63	127	120.139	6.8612	47.076
2	64	121	126.276	-5.2764	27.840
3	66	142	138.552	3.4484	11.891
4	69	157	156.964	0.0356	0.001
5	69	162	156.964	5.0356	25.357
6	71	156	169.240	-13.2396	175.287
7	71	169	169.240	-0.2396	0.057
8	72	165	175.377	-10.3772	107.686
9	73	181	181.515	-0.5148	0.265
10	75	208	193.790	14.2100	201.924
					597.4

$$e_i = y_i - \hat{y}_i$$

$$Q = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Least Squares

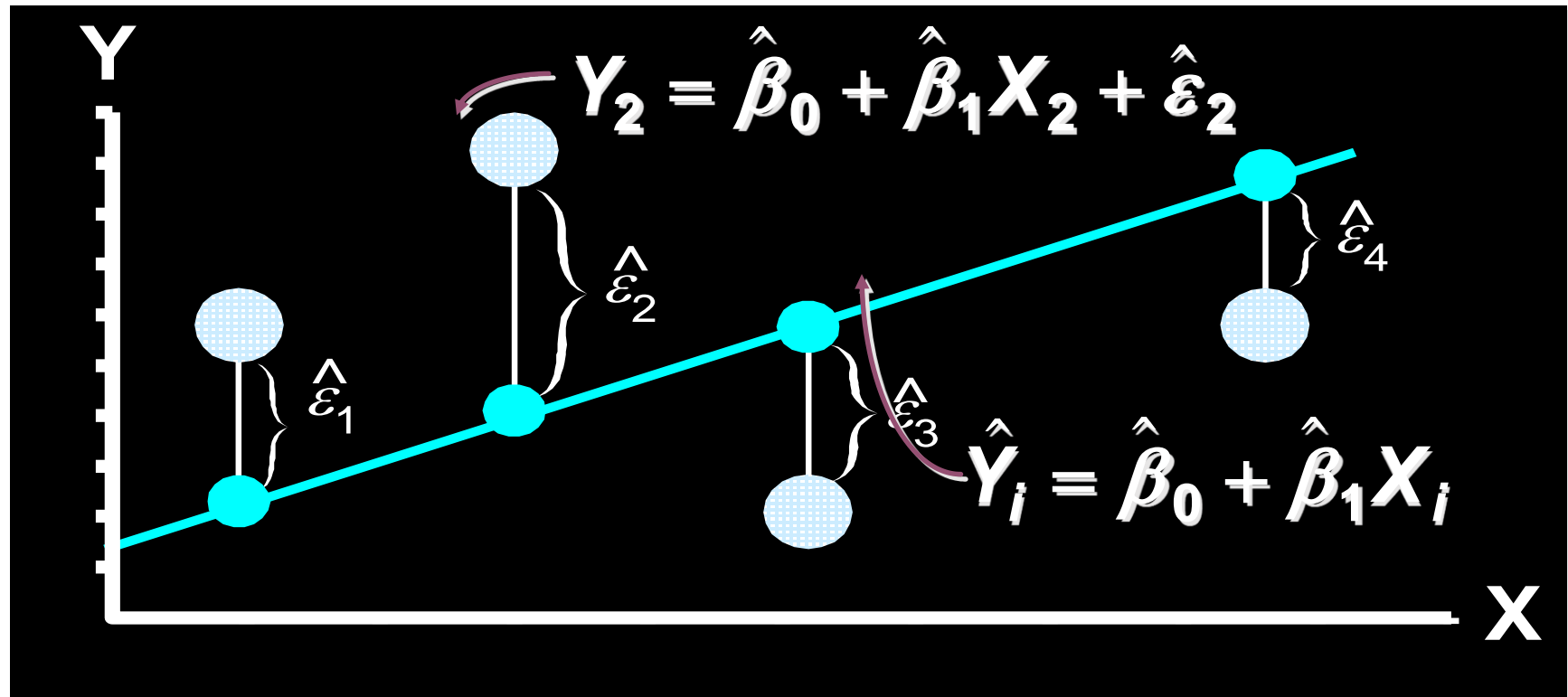
- 1. 'Best Fit' Means Difference Between Actual Y Values & Predicted Y Values Are a Minimum. *But* Positive Differences Off-Set Negative. So square errors!

$$\sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n \hat{\epsilon}_i^2$$

- 2. LS Minimizes the Sum of the Squared Differences (errors) (SSE)

Least Squares Graphically

$$\text{LS minimizes } \sum_{i=1}^n \hat{\varepsilon}_i^2 = \hat{\varepsilon}_1^2 + \hat{\varepsilon}_2^2 + \hat{\varepsilon}_3^2 + \hat{\varepsilon}_4^2$$



How to estimate parameters

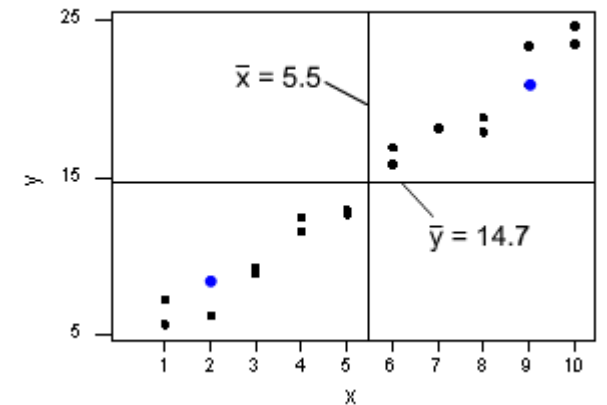
We minimize the equation for the sum of the squared prediction errors:

$$Q = \sum_{i=1}^n (y_i - (b_0 + b_1 x_i))^2$$

(that is, take the derivative with respect to b_0 and b_1 , set to 0, and solve for b_0 and b_1) and get the "least squares estimates" for b_0 and b_1 :

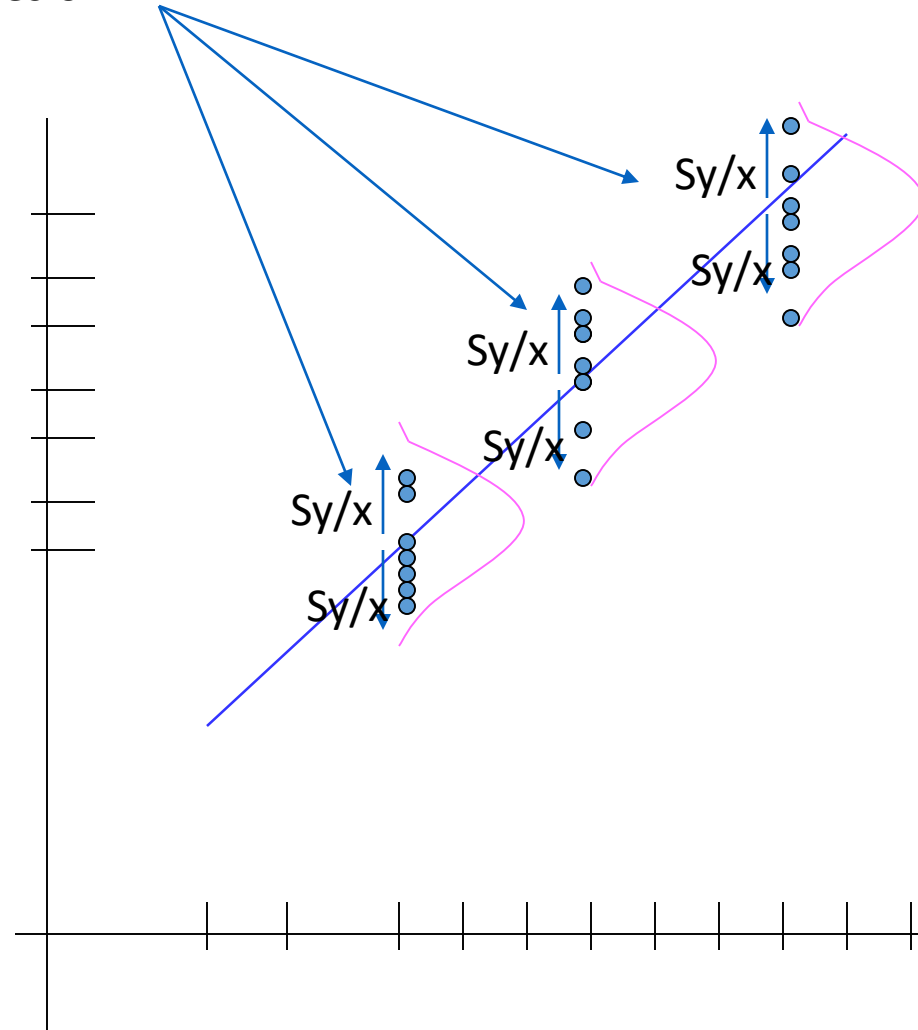
$$b_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

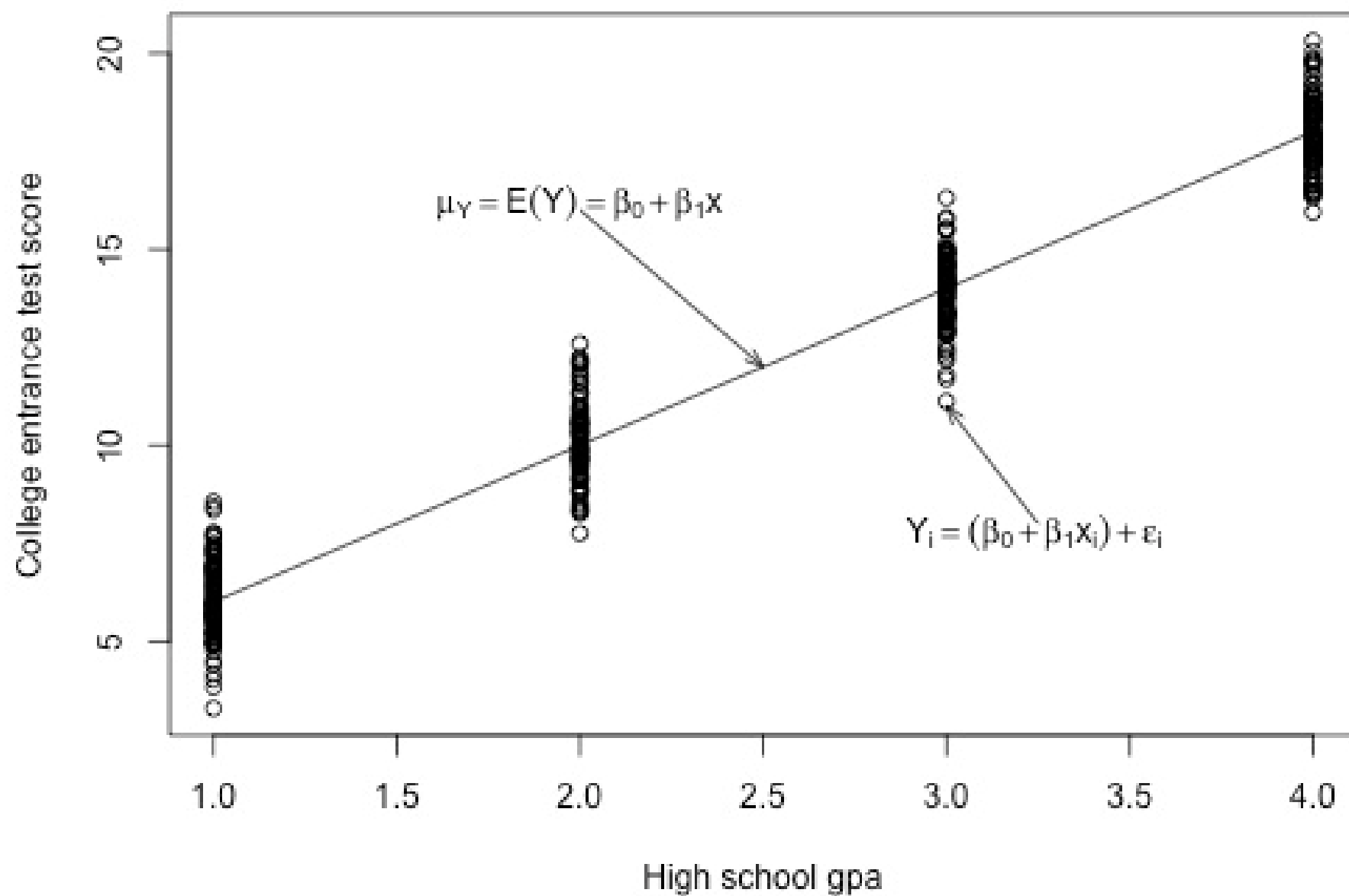
$$b_0 = \bar{y} - b_1 \bar{x}$$



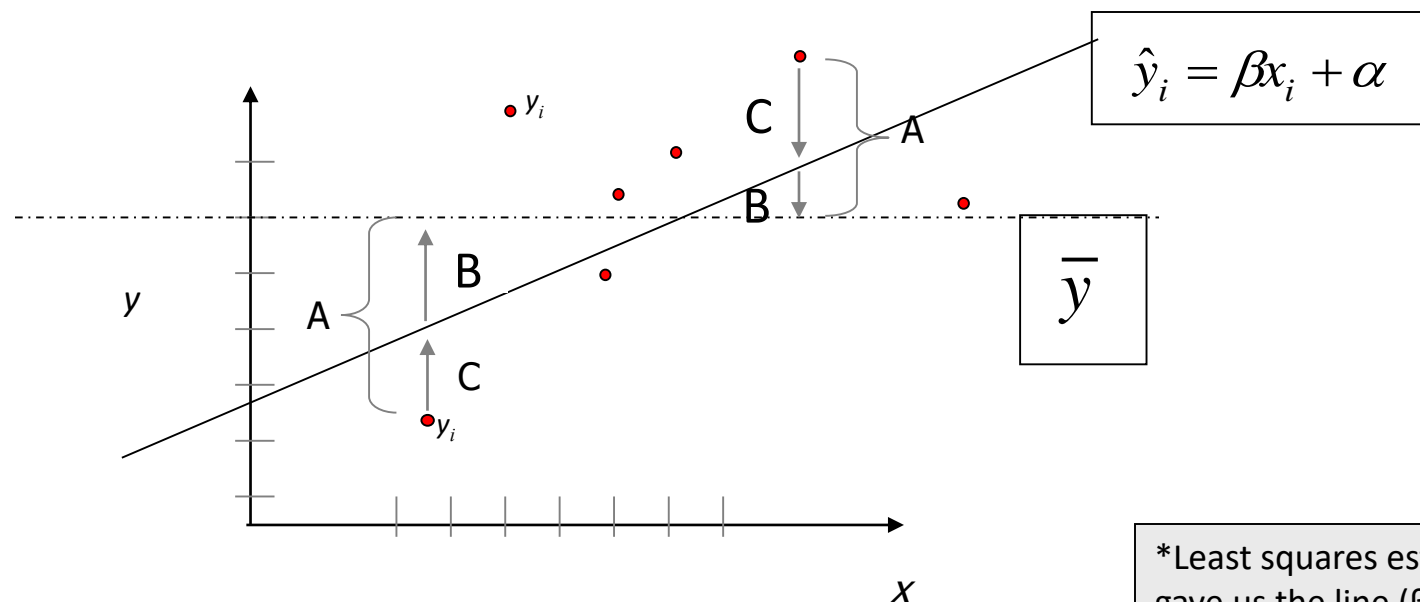
the least squares line passes through the point (\bar{x}, \bar{y}) , since when $x = \bar{x}$, then $y = b_0 + b_1 \bar{x} = \bar{y} - b_1 \bar{x} + b_1 \bar{x} = \bar{y}$.

The standard error of Y given X is the average variability around the regression line at any given value of X. It is assumed to be equal at all values of X.





Regression Picture



*Least squares estimation gave us the line (β) that minimized C^2

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

A^2 B^2 C^2

$$R^2 = SS_{\text{reg}} / SS_{\text{total}}$$

SS_{total}
Total squared distance of observations from naïve mean of y
Total variation

SS_{reg}
Distance from regression line to naïve mean of y
Variability due to x (regression)

SS_{residual}
Variance around the regression line
Additional variability not explained by x—what least squares method aims to minimize

Regression Line

- If the scatter plot of our sample data suggests a linear relationship between two variables i.e.

$$y = \beta_0 + \beta_1 x$$

we can summarize the relationship by drawing a straight line on the plot.

- Least squares method give us the “best” estimated line for our set of sample data.

Regression Line

- We will write an estimated regression line based on sample data as

$$\hat{y} = b_0 + b_1x$$

- The method of least squares chooses the values for b_0 , and b_1 to minimize the sum of squared errors

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y - b_0 - b_1x)^2$$

Regression Line

- Using calculus, we obtain estimating formulas:

or

$$b_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

$$b_1 = r \frac{S_y}{S_x}$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

Estimation of Mean Response

- Fitted regression line can be used to estimate the mean value of y for a given value of x .
- Example
 - The weekly advertising expenditure (x) and weekly sales (y) are presented in the following table.

y	x
1250	41
1380	54
1425	63
1425	54
1450	48
1300	46
1400	62
1510	61
1575	64
1650	71

Point Estimation of Mean Response

- From previous table we have:

$$\begin{aligned} n &= 10 & \sum x &= 564 & \sum x^2 &= 32604 \\ \sum y &= 14365 & \sum xy &= 818755 \end{aligned}$$

- The least squares estimates of the regression coefficients are:

$$b_1 = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2} = \frac{10(818755) - (564)(14365)}{10(32604) - (564)^2} = 10.8$$

$$b_0 = 1436.5 - 10.8(56.4) = 828$$

Point Estimation of Mean Response

- The estimated regression function is:

$$\hat{y} = 828 + 10.8x$$

$$\text{Sales} = 828 + 10.8 \text{ Expenditure}$$

- This means that if the weekly advertising expenditure is increased by \$1 we would expect the weekly sales to increase by \$10.8.

Point Estimation of Mean Response

- Fitted values for the sample data are obtained by substituting the x value into the estimated regression function.
- For example if the advertising expenditure is \$50, then the estimated Sales is:

$$Sales = 828 + 10.8(50) = 1368$$

- This is called the point estimate (forecast) of the mean response (sales).

Linear correlation and linear regression

Covariance

$$\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y})}{n - 1}$$

Interpreting Covariance

$\text{cov}(X,Y) > 0$ ~~X~~→ and Y are positively correlated

$\text{cov}(X,Y) < 0$ ←~~X~~ and Y are inversely correlated

$\text{cov}(X,Y) = 0$ ~~X~~→ and Y are independent

Correlation coefficient

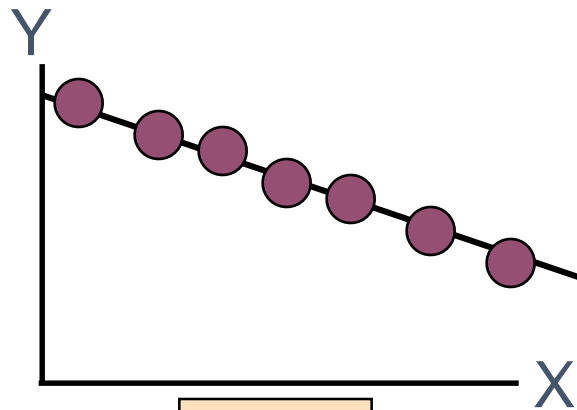
- Pearson's Correlation Coefficient is standardized covariance (unitless):

$$r = \frac{\text{covariance}(x, y)}{\sqrt{\text{var } x} \sqrt{\text{var } y}}$$

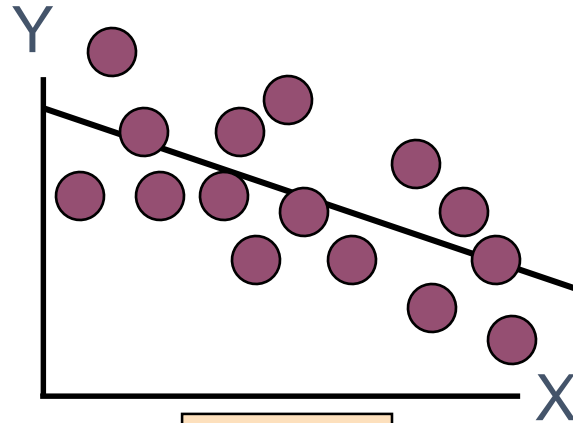
Correlation

- Measures the relative strength of the *linear* relationship between two variables
- Unit-less
- Ranges between -1 and 1
- The closer to -1 , the stronger the negative linear relationship
- The closer to 1 , the stronger the positive linear relationship
- The closer to 0 , the weaker any positive linear relationship

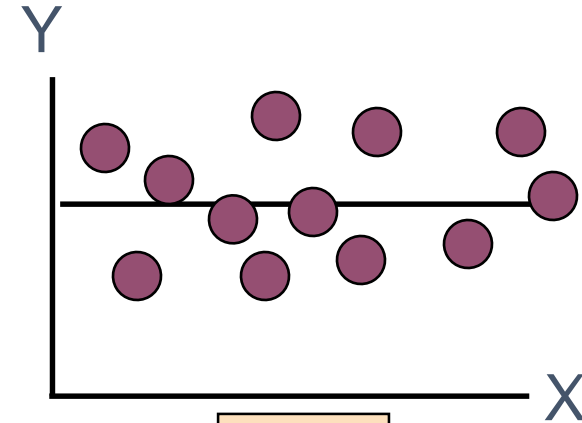
Scatter Plots of Data with Various Correlation Coefficients



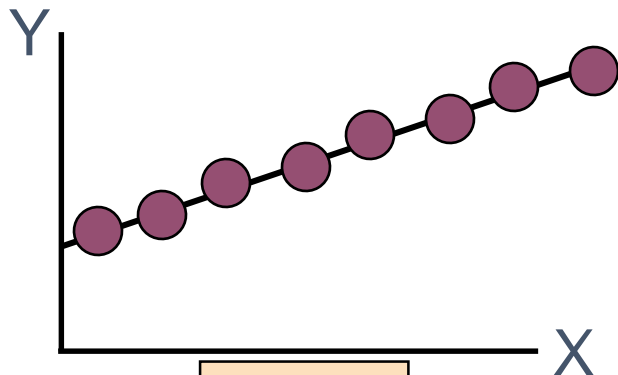
$$r = -1$$



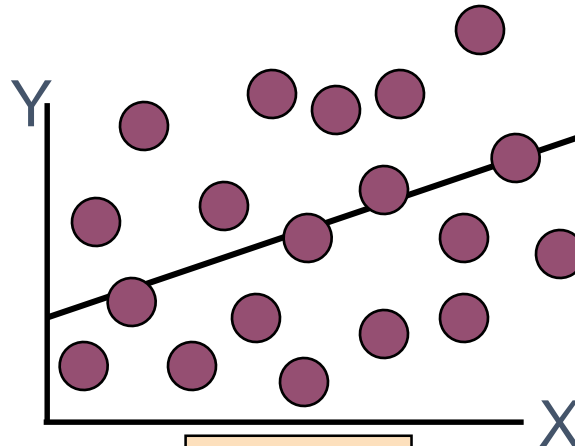
$$r = -.6$$



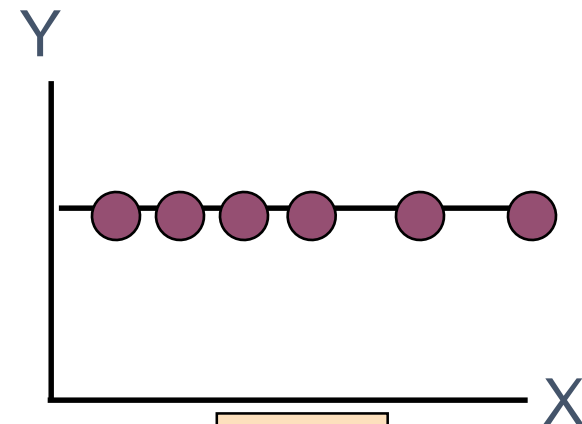
$$r = 0$$



$$r = +1$$



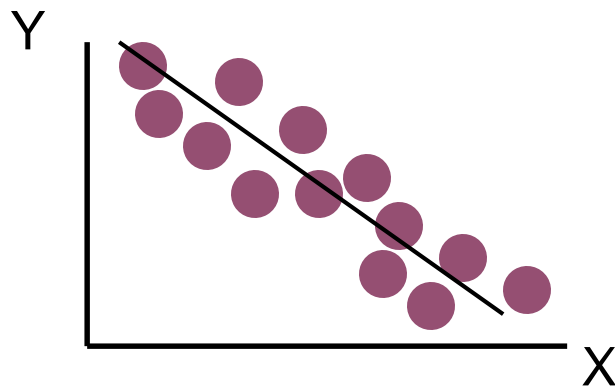
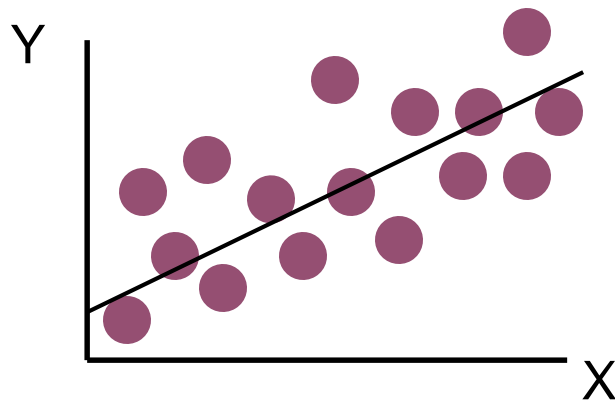
$$r = +.3$$



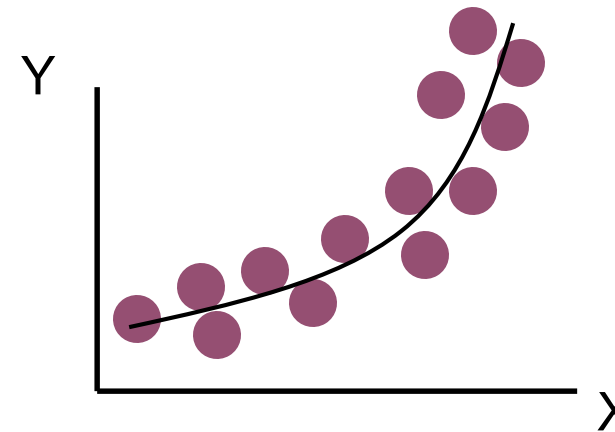
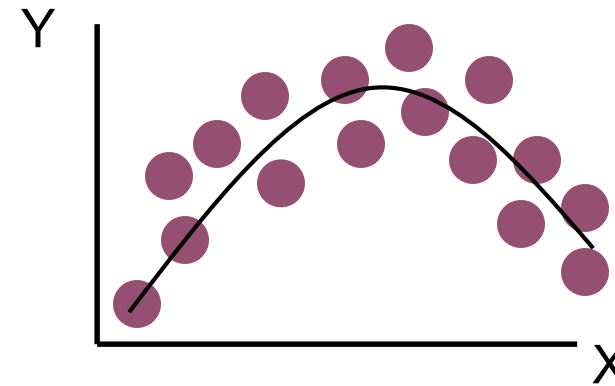
$$r = 0$$

Linear Correlation

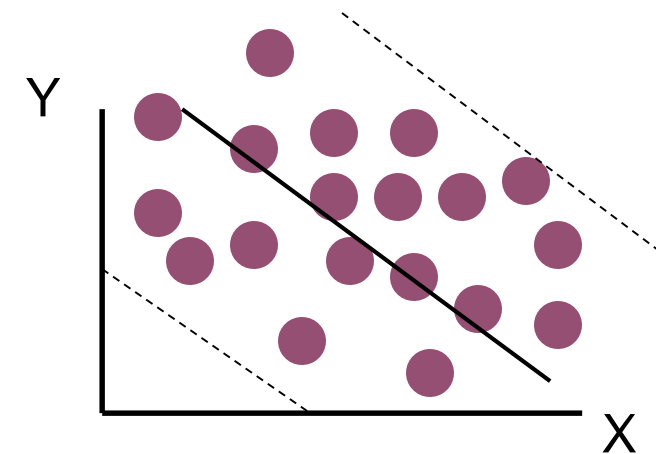
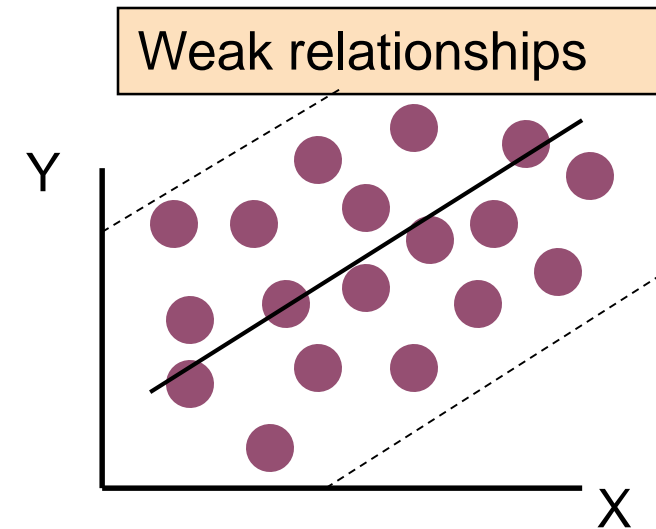
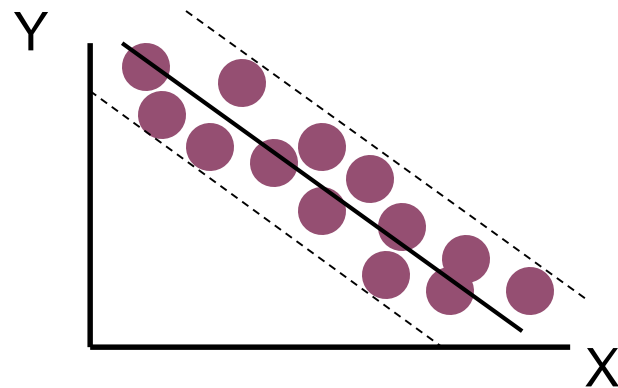
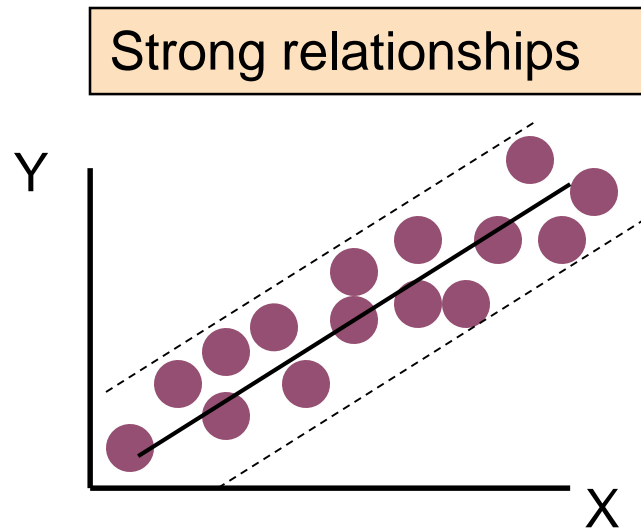
Linear relationships



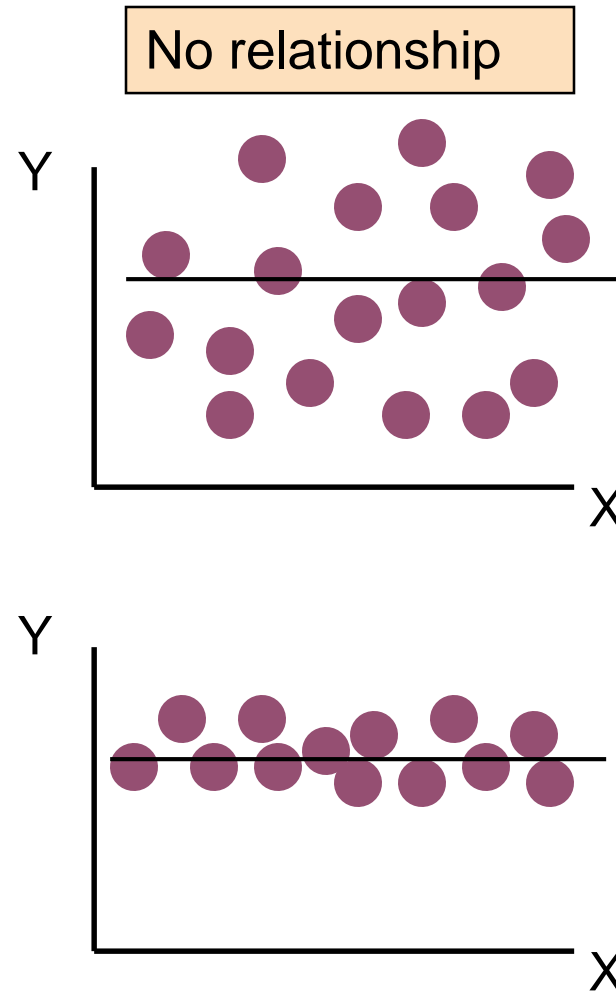
Curvilinear relationships



Linear Correlation



Linear Correlation



Calculating by hand...

$$\hat{r} = \frac{\text{covariance}(x, y)}{\sqrt{\text{var } x} \sqrt{\text{var } y}} = \frac{\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}}{\sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}} \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}}}$$

Simpler calculation formula...

$$\hat{r} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} = \frac{SS_{xy}}{\sqrt{SS_x SS_y}}$$

Numerator of
covariance

$$\hat{r} = \frac{SS_{xy}}{\sqrt{SS_x SS_y}}$$

Numerators of
variance

Least Square estimation

Slope (beta coefficient) =

$$\hat{\beta} = \frac{Cov(x, y)}{Var(x)}$$

Intercept=

$$\text{Calculate : } \hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}$$

Regression line always goes through the point: (\bar{x}, \bar{y})

Relationship with correlation

$$\hat{r} = \hat{\beta} \frac{SD_x}{SD_y}$$

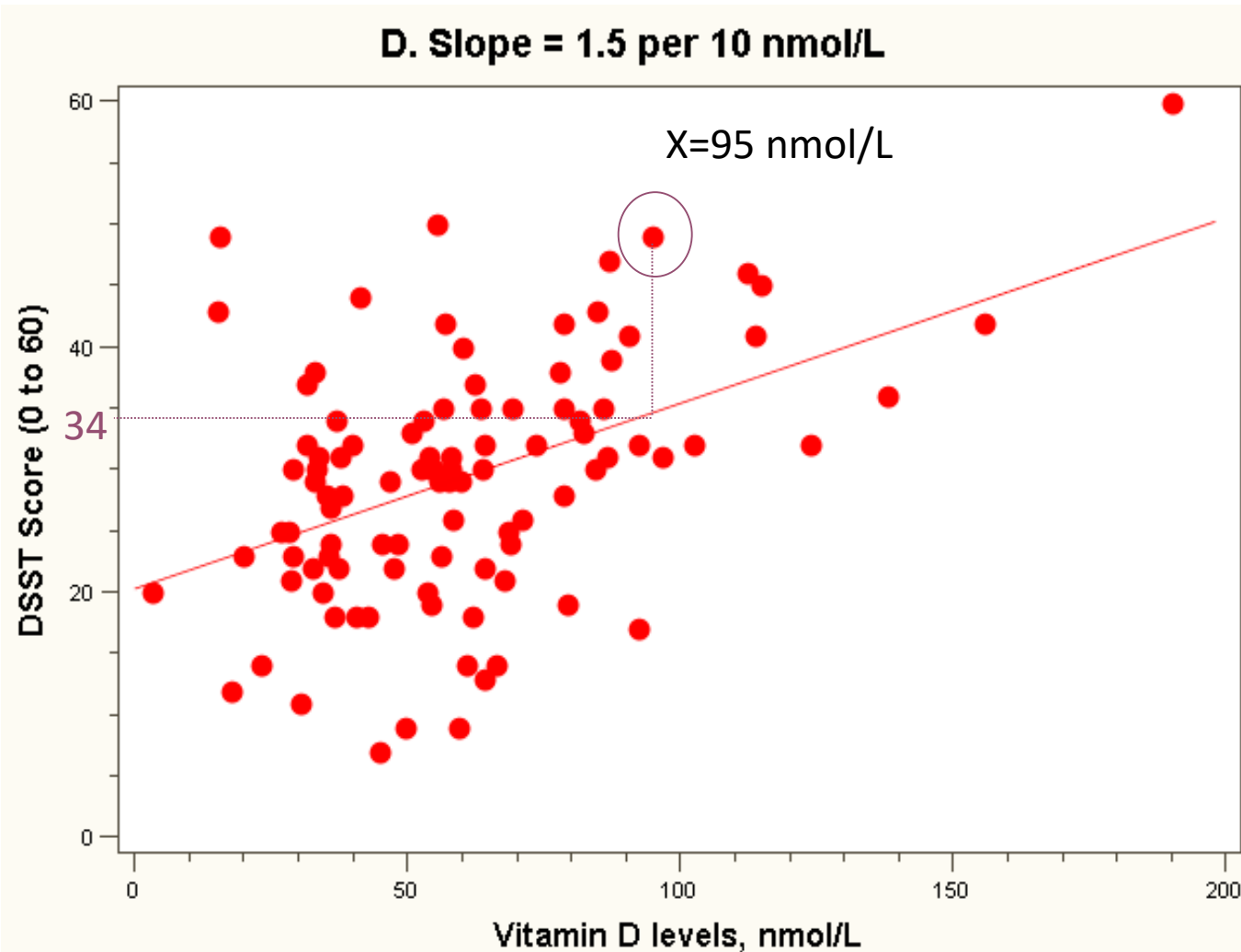
In correlation, the two variables are treated as equals. In regression, one variable is considered independent (=predictor) variable (X) and the other the dependent (=outcome) variable Y.

Residual Analysis: check assumptions

$$e_i = Y_i - \hat{Y}_i$$

- The residual for observation i , e_i , is the difference between its observed and predicted value
- Residuals are highly useful for studying whether a given regression model is appropriate for the data at hand.
- Check the assumptions of regression by examining the residuals
 - Examine for linearity assumption
 - Examine for constant variance for all levels of X (homoscedasticity)
 - Evaluate normal distribution assumption
 - Evaluate independence assumption
- Graphical Analysis of Residuals

Residual =
observed - predicted

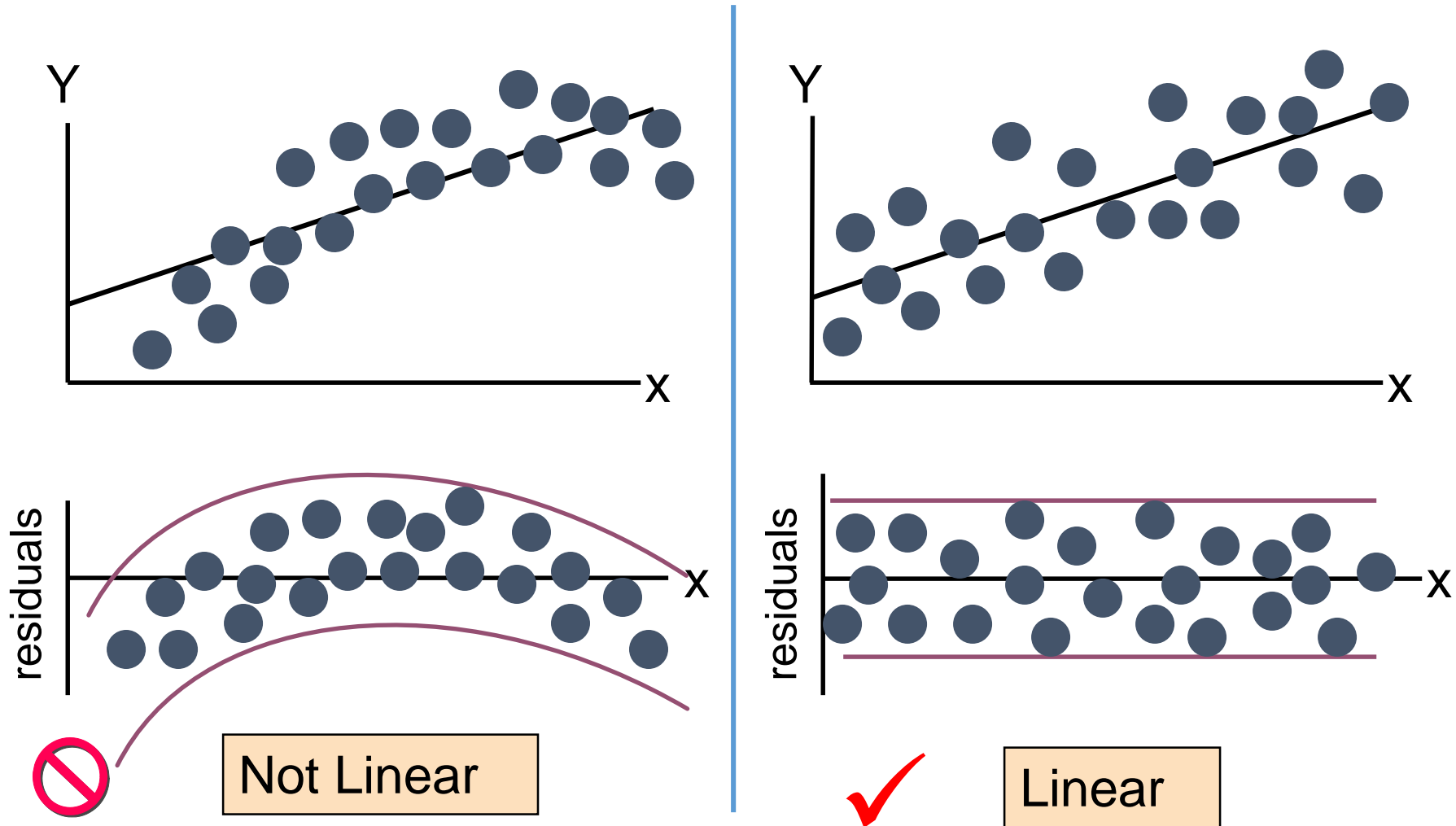


$$y_i = 48$$

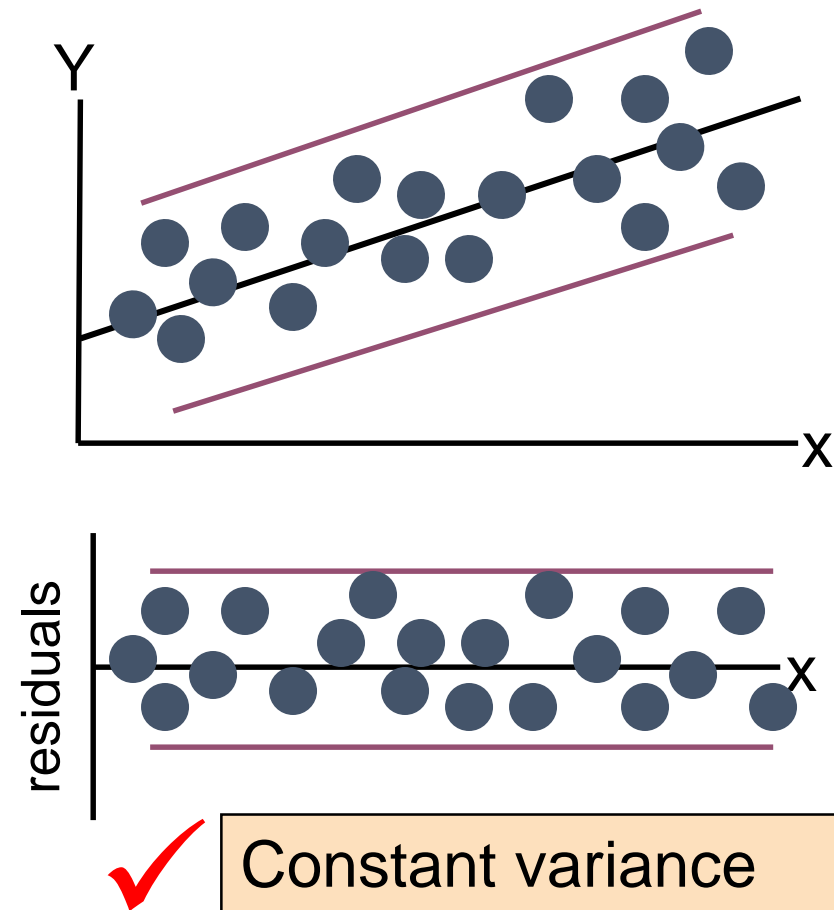
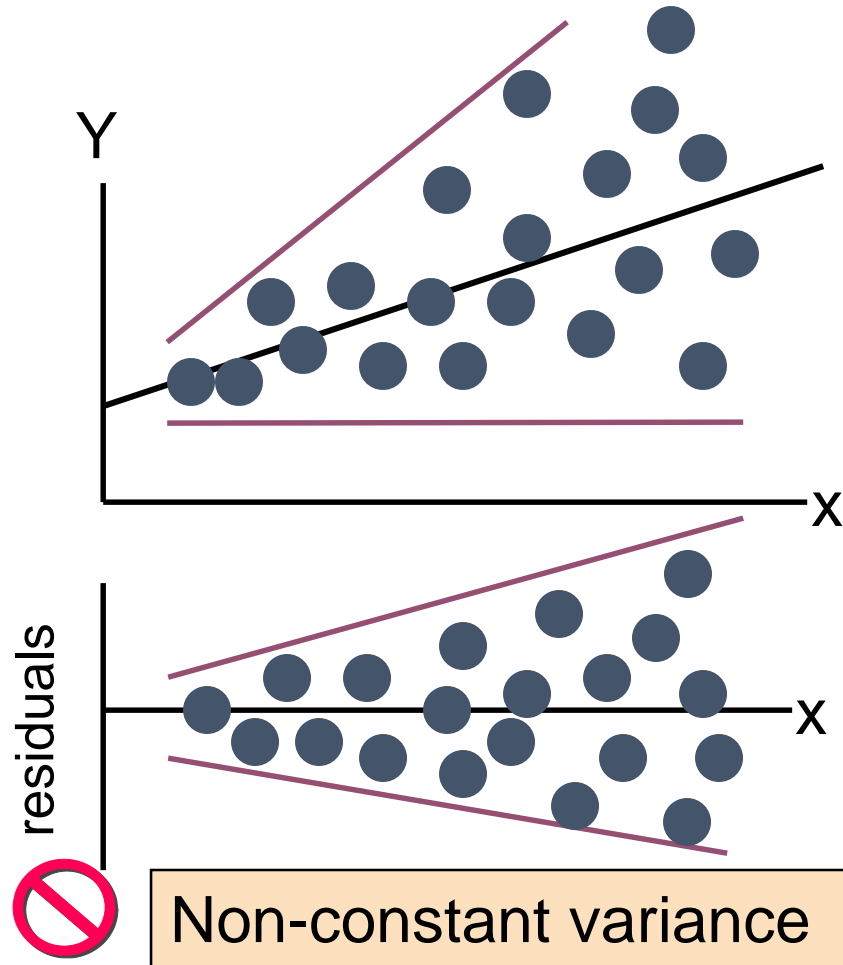
$$\hat{y}_i = 34$$

$$y_i - \hat{y}_i = 14$$

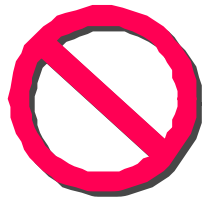
Residual Analysis for Linearity



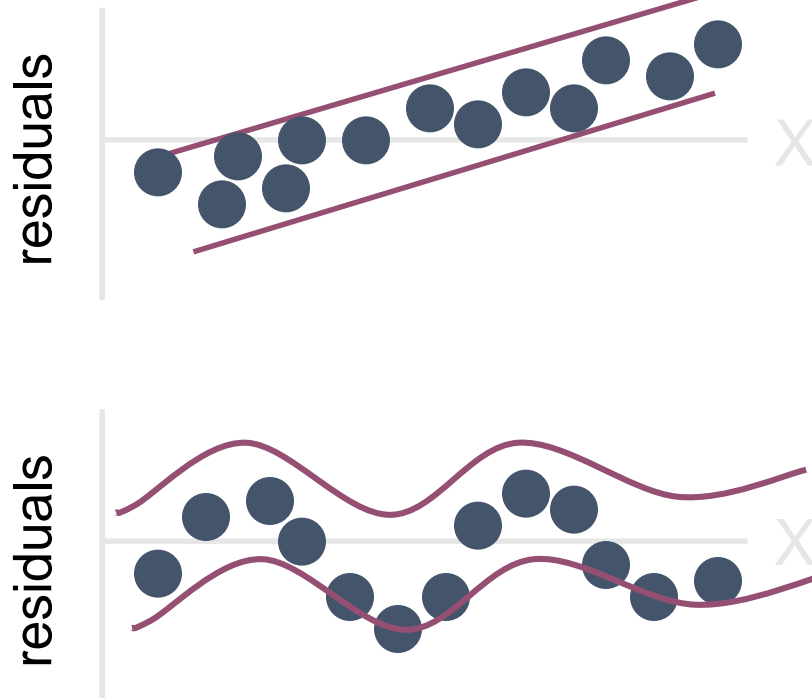
Residual Analysis for Homoscedasticity



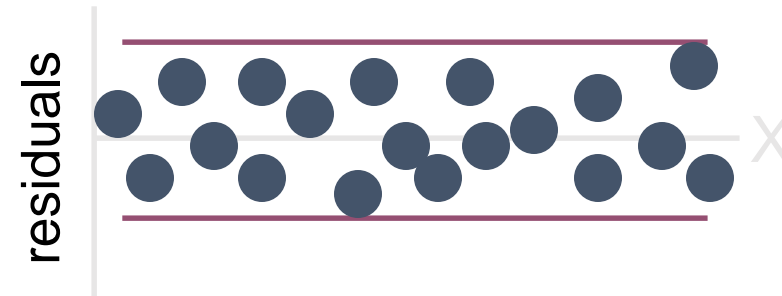
Residual Analysis for Independence



Not Independent



Independent



Example: weekly advertising expenditure

y	x	y-hat	Residual (e)
1250	41	1270.8	-20.8
1380	54	1411.2	-31.2
1425	63	1508.4	-83.4
1425	54	1411.2	13.8
1450	48	1346.4	103.6
1300	46	1324.8	-24.8
1400	62	1497.6	-97.6
1510	61	1486.8	23.2
1575	64	1519.2	55.8
1650	71	1594.8	55.2

Estimation of the variance of the error terms, σ^2

- The variance σ^2 of the error terms ε_i in the regression model needs to be estimated for a variety of purposes.
 - It gives an indication of the variability of the probability distributions of y .
 - It is needed for making inference concerning regression function and the prediction of y .

Regression Standard Error

- To estimate σ we work with the variance and take the square root to obtain the standard deviation.
- For simple linear regression the estimate of σ^2 is the average squared residual.

$$s_{y.x}^2 = \frac{1}{n-2} \sum e_i^2 = \frac{1}{n-2} \sum (y_i - \hat{y}_i)^2$$

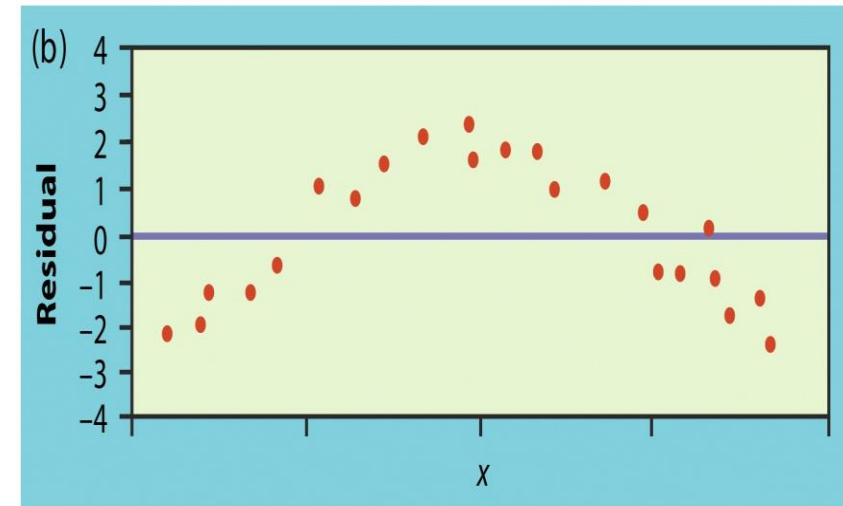
- To estimate σ , use
- s estimates the standard deviation σ of the error term ε in the statistical model for simple linear regression. $s_{y.x} = \sqrt{s_{y.x}^2}$

Regression Standard Error

y	x	y-hat	Residual (e)	square(e)
1250	41	1270.8	-20.8	432.64
1380	54	1411.2	-31.2	973.44
1425	63	1508.4	-83.4	6955.56
1425	54	1411.2	13.8	190.44
1450	48	1346.4	103.6	10732.96
1300	46	1324.8	-24.8	615.04
1400	62	1497.6	-97.6	9525.76
1510	61	1486.8	23.2	538.24
1575	64	1519.2	55.8	3113.64
1650	71	1594.8	55.2	3047.04
y-hat = 828+10.8X			total	36124.76
			S _{y .x}	67.19818

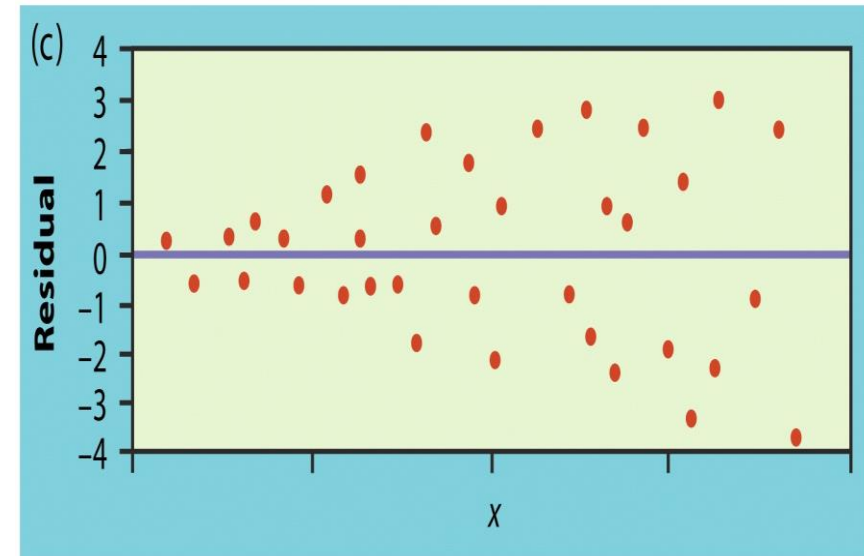
Residual plots

- The points in this residual plot have a curve pattern, so a straight line fits poorly



Residual plots

- The points in this plot show more spread for larger values of the explanatory variable x , so prediction will be less accurate when x is large.



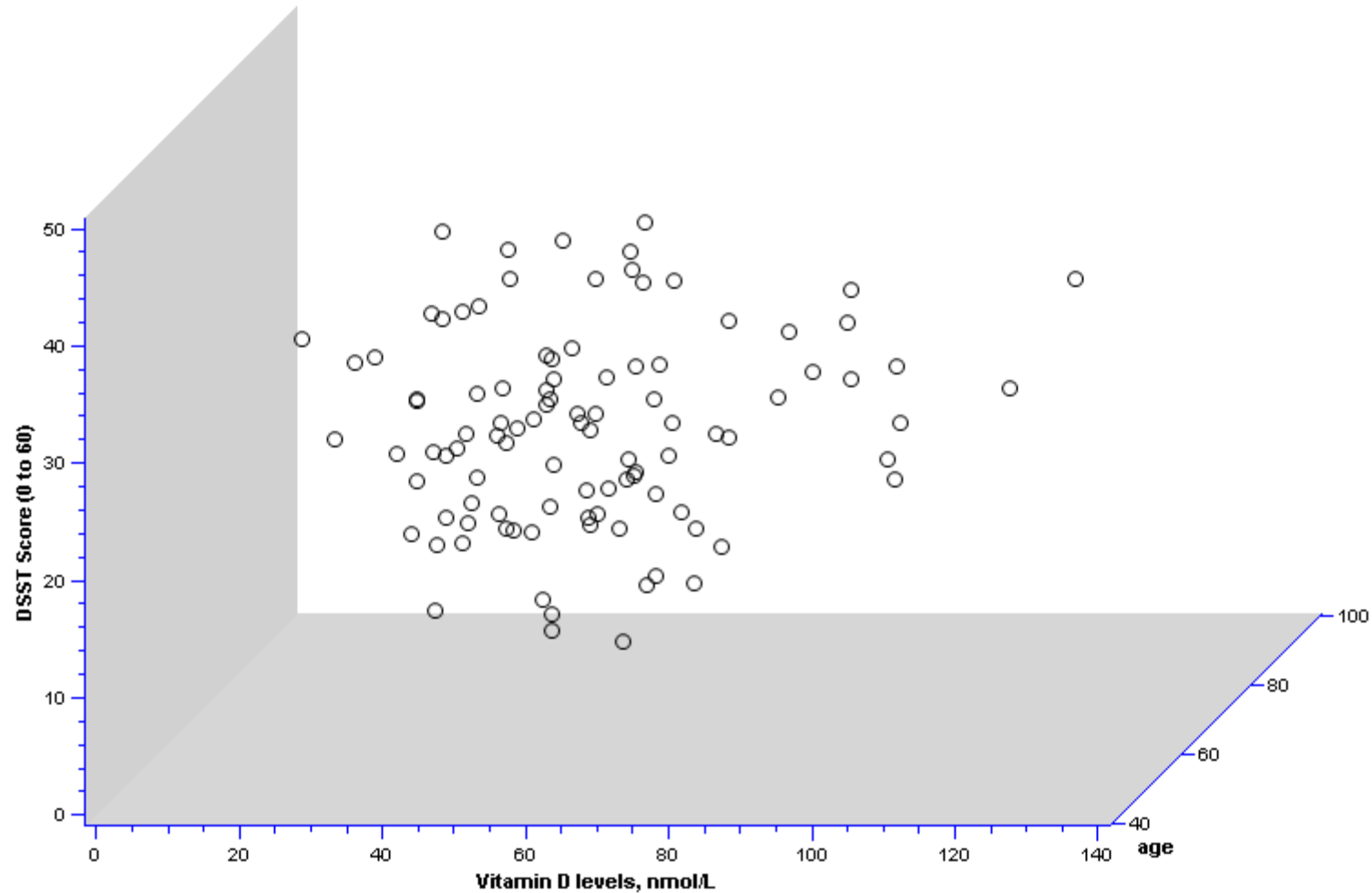
Variable transformations

- If the residual plot suggests that the variance is not constant, a transformation can be used to stabilize the variance.
- If the residual plot suggests a non linear relationship between x and y, a transformation may reduce it to one that is approximately linear.
- Common linearizing transformations are:
- Variance stabilizing transformations are:

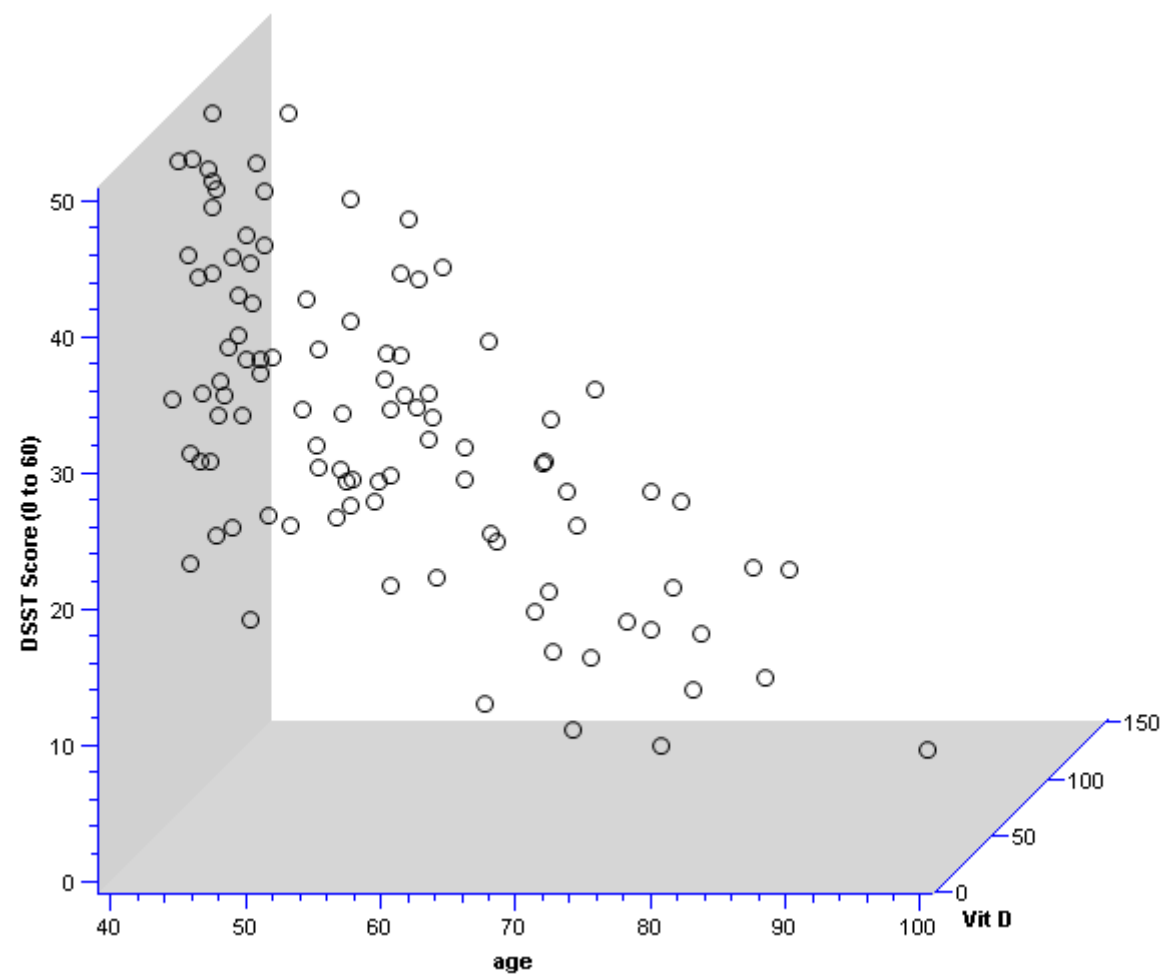
$$\frac{1}{x}, \quad \log(x)$$

$$\frac{1}{y}, \quad \log(y), \quad \sqrt{y}, \quad y^2$$

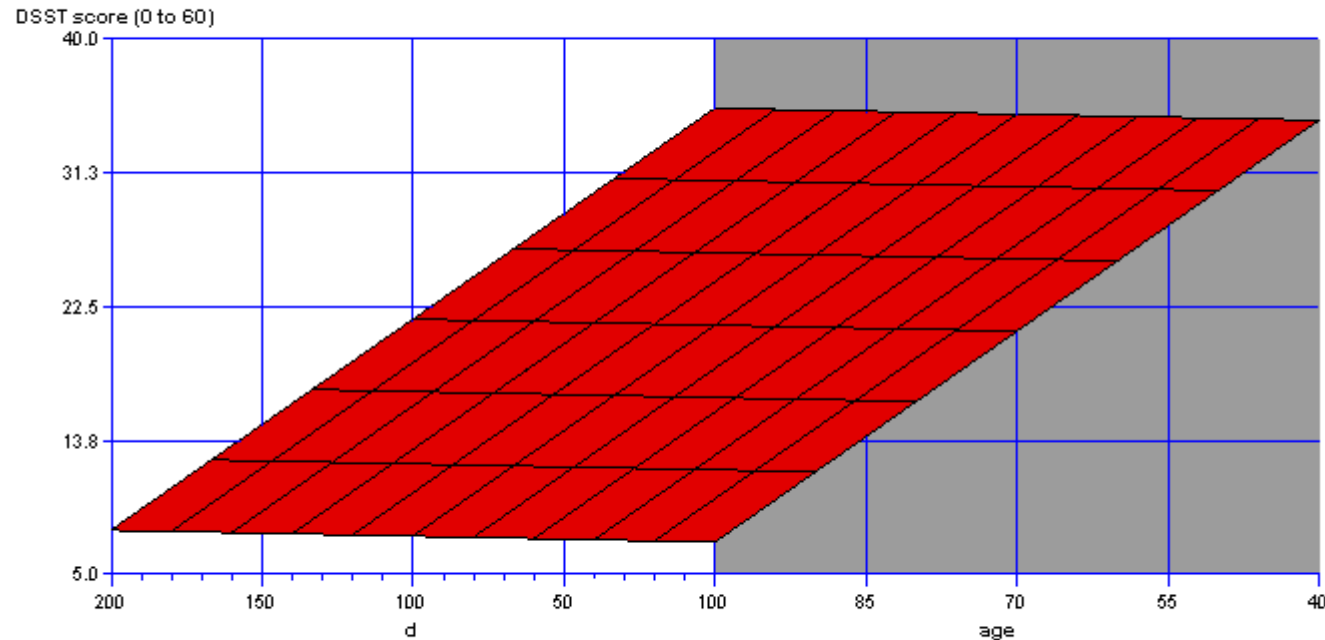
2 predictors: age and vit D...



Different 3D view...



Fit a plane rather than a line...



On the plane, the slope for vitamin D is the same at every age; thus, the slope for vitamin D represents the effect of vitamin D when age is held constant.